

THE EQUATIONAL THEORIES PROJECT: ADVANCING COLLABORATIVE MATHEMATICAL RESEARCH AT SCALE

MATTHEW BOLAN, JOACHIM BREITNER, JOSE BROX, MARIO CARNEIRO, MARTIN DVOŘÁK, ANDRÉS GOENS, AARON HILL, HARALD HUSUM, ZOLTAN KOCSIS, BRUNO LE FLOCH, LORENZO LUCCIOLI, ALEX MEIBURG, PIETRO MONTICONE, GIOVANNI PAOLINI, BERNHARD REINKE, DAVID RENSHAW, MARCUS ROSSEL, CODY ROUX, JÉRÉMY SCANVIC, SHREYAS SRINIVAS, ANAND RAO TADIPATRI, TERENCE TAO, VLAD TSYRKLEVICH, DANIEL WEBER, FAN ZHENG

ABSTRACT. We report on the *Equational Theories Project* (ETP), an online collaborative pilot project to explore new ways to collaborate in mathematics with machine assistance. The project successfully determined all 22 028 942 edges of the implication graph between the 4694 simplest equational laws on magmas, by a combination of human-generated and automated proofs, all validated by the formal proof assistant language *Lean*. As a result of this project, several new constructions of magmas obeying specific laws were discovered, and several auxiliary questions were also addressed, such as the effect of restricting attention to finite magmas.

CONTENTS

| | |
|--|----|
| 1. Introduction | 3 |
| 1.1. Magmas and Equational Laws | 3 |
| 1.2. The Equational Theories Project | 5 |
| 1.3. Outcomes | 6 |
| 1.4. Further directions | 9 |
| 2. Notation and Mathematical Foundations | 9 |
| 3. Formal Foundations | 11 |
| 4. Project Management | 13 |
| 4.1. Problems of scale in mathematical collaboration | 13 |
| 4.2. Handling Scaling Issues | 14 |
| 4.3. Other Design Considerations | 15 |

Date: May 23, 2025.

| | |
|---|----|
| 4.4. Maintenance | 15 |
| 5. Counterexample constructions | 16 |
| 5.1. Finite magmas | 16 |
| 5.2. Linear models | 17 |
| 5.3. Translation-invariant models | 19 |
| 5.4. The twisting semigroup | 20 |
| 5.5. Greedy constructions | 22 |
| 5.6. Modifying base models | 26 |
| 6. Syntactic arguments | 28 |
| 6.1. Simple rewrites | 28 |
| 6.2. Matching invariants | 28 |
| 6.3. Confluence | 30 |
| 6.4. Unique factorization | 33 |
| 7. Proof Automation | 36 |
| 7.1. Proof Techniques | 37 |
| 7.2. Proof Reconstruction and Integration | 38 |
| 7.3. Empirical Results | 40 |
| 8. Implications for Finite Magmas | 41 |
| 9. Higman–Neumann laws | 41 |
| 10. AI and Machine Learning Contributions | 41 |
| 10.1. Graph ML: Directed link prediction on the implication graph | 41 |
| 11. User Interfaces | 47 |
| 12. Data Management | 48 |
| 13. Conclusions and Future Directions | 49 |
| 13.1. Miscellaneous remarks | 52 |

| | |
|----------------------------------|----|
| Acknowledgments | 52 |
| Appendix A. Numbering system | 52 |
| Appendix B. Author Contributions | 55 |
| References | 56 |
| Todo list | 58 |

1. INTRODUCTION

The purpose of this paper is to report on the *Equational Theories Project* (ETP)¹, a pilot project launched² in September 2024 to explore new ways to collaboratively work on mathematical research projects using machine assistance. The project goal, in the area of universal algebra, was selected³ to be particularly amenable to crowdsourced and computer-assisted techniques, while still being of mathematical research interest.

The project achieved its primary goal on 14 April 2025, when the $4694 \times (4694 - 1) = 22\,028\,942$ implications between the test set of 4964 equational laws were completely determined, with proofs or refutations formalized in *Lean*. This required coordinating the efforts of a large number of participants contributing both human-written formalizations and automatically generated proofs from various computer tools. In this paper, we report on both the scientific outcomes of the project, as well as the organizational issues that came up with organizing a mathematical project of this scale.

1.1. Magmas and Equational Laws. In order to describe the mathematical goals of the ETP, we need some notation. A *magma* $\mathcal{M} = (M, \diamond)$ is a set M (known as the *carrier*) together with a binary operation $\diamond: M \times M \rightarrow M$. An *equational law* for a magma, or *law* for short, is an identity involving \diamond and some formal indeterminates, which we will typically denote using the Roman letters x, y, z, w, u, v , as well as the formal equality symbol \simeq in place of the equality symbol $=$ to emphasize the formal nature of the law.

In the ETP, a unique number was assigned to each equational law, via a numbering system that we describe in Section A. For instance, the *commutative law* $x \diamond y \simeq y \diamond x$ is assigned the equation number (E43), while the *associative law* $(x \diamond y) \diamond z \simeq x \diamond (y \diamond z)$ is assigned the equation number (E4512). A list of all equations referred to by number in this paper is also provided in Section A.

A magma $\mathcal{M} = (M, \diamond)$ obeys a law E if the law E holds for all possible assignments of the indeterminate to elements of M , in which case we write $\mathcal{M} \models E$. Thus, for instance

¹https://teorth.github.io/equational_theories/

²<https://terrytao.wordpress.com/2024/09/25>

³The specific mathematical goal was inspired by a MathOverflow question.

$\mathcal{M} \models \text{E43}$ if one has $x \diamond y = y \diamond x$ for all $x, y \in M$. Note that the formal indeterminate symbols x, y in E43 are now replaced by concrete elements x, y of the carrier M .

We say that a law E *entails* or *implies* another law E' if every magma that obeys E , also implies E' : $(\mathcal{M} \models E) \implies (\mathcal{M} \models E')$. We write this relation as $E \vdash E'$. We say that two laws are *equivalent* if they entail each other. For instance, the constant law $x \diamond y \simeq z \diamond w$ (E46) can easily be seen to be equivalent to the law $x \diamond x \simeq y \diamond z$ (E41). It is clear that \vdash is a pre-order, that is to say a partial order after one quotients by equivalence.

In this entailment pre-ordering, the maximal element is given by the trivial law $x \simeq x$ (E1), and the minimal element is given by the singleton law $x \simeq y$ (E2), thus $\text{E2} \vdash E \vdash \text{E1}$ for all laws E .

We also define a variant: we say that E *entails E' for finite magmas*, and write $E \vdash_{\text{fin}} E'$, if every *finite* magma M that obeys E , also obeys E' . Clearly, the relation $E \vdash E'$ implies $E \vdash_{\text{fin}} E'$; but, as observed by Austin [2], the converse is not true in general.

The *order* of an equational law is the number of occurrences of the magma operation, and can be viewed as a crude measure of complexity of the law. For instance, the commutative law (E43) has order 2, while the associative law (E4512) has order 4. We note some selected laws of small order that have previously appeared in the literature:

- The *central groupoid law* $x \simeq (y \diamond x) \diamond (x \diamond z)$ (E168) is an order-3 law introduced by Evans [11] and studied further by Knuth [19] and many further authors, being closely related to central digraphs (also known as unique path property digraphs), and leading in particular to the discovery of the Knuth-Bendix algorithm [20]; see [25] for a more recent survey.
- *Tarski's axiom* $x \simeq y \diamond ((z \diamond (x \diamond (y \diamond z))))$ (E543) is an order-4 law that was shown by Tarski [42] to characterize the operation of subtraction in an abelian group; that is to say, a magma $\mathcal{M} = (M, \diamond)$ obeys (E543) if and only if there is an abelian group structure on \mathcal{M} for which $x \diamond y = x - y$ for all $x, y \in M$.
- In a similar vein, it was shown in [33] (see also [34]) that the order-4 law $x \simeq (y \diamond z) \diamond (y \diamond (x \diamond z))$ (E1571) characterizes addition (or subtraction) in an abelian group of exponent 2; it was shown in [31] that the order-6 law $x \simeq (y \diamond ((x \diamond y) \diamond y)) \diamond (x \diamond (z \diamond y))$ (E345169) characterizes the Sheffer stroke in a boolean algebra, and it was shown in [14] that the order-8 law $x \simeq y \diamond (((y \diamond y) \diamond x) \diamond z) \diamond (((y \diamond y) \diamond y) \diamond z)$ (E42323216) characterizes division in a (not necessarily abelian) group.

Some further examples of laws characterizing well-known algebraic structures are listed in [30].

The Birkhoff completeness theorem [3, Th. 3.5.14] implies that an implication $E \vdash E'$ of equational laws holds if and only if the left-hand side of E' can be transformed into the right-hand side by a finite number of substitution rewrites using the law E . However, the problem of determining whether such an implication holds is undecidable in general [32]. Even when the order is small, some implications⁴ can require lengthy computer-assisted

⁴Another contemporaneous example of this phenomenon was the solution of the Robbins problem [29].

proofs; for instance, it was noted in [16] that the order-4 law $x \simeq (y \diamond x) \diamond ((x \diamond z) \diamond z)$ (E1689) was equivalent to the singleton law (E2), but all known proofs were found with computer assistance.⁵ Furthermore, for the finite magma implication relation $E \vdash_{\text{fin}} E'$, no analogue of the Birkhoff completeness theorem is available.

1.2. The Equational Theories Project. As noted in Section A, there are 4694 equational laws of order at most 4. The primary mathematical goal of the ETP was to completely determine the *implication graph* for these laws, in which there is a directed edge from E to E' precisely when $E \vdash E'$. As the project progressed, an additional goal was added to determine the slightly larger *finite implication graph*, in which there is a directed edge from E to E' precisely when $E \vdash_{\text{fin}} E'$.

Such systematic determinations of implication graphs have been seen previously in the literature; for instance, in [37], the relations between 60 identities of Bol–Moufang type were established, and in the blog post [47, §17], some initial steps towards generating this graph for the first hundred or so laws on our list were performed. However, to our knowledge, the ETP is the first project to study such implications at the scale of thousands of laws.

The ETP requires the determination of the truth or falsity of $4694^2 = 22\,033\,636$ implications (for both arbitrary magmas and finite magmas), or $4694 \times (4694 - 1) = 22\,028\,942$ if the reflexive implications $E \vdash E$ are removed; while one can use properties such as the transitivity of entailment to reduce the work somewhat, this is clearly a task that requires significant automation. It was also a project highly amenable to crowdsourcing, in which different participants could work on developing different techniques, each of which could be used to fill out a different part of the implication graph. In this respect, the project could be compared with a Polymath project [13], which used online forums such as blogs and wikis to openly collaborate on a mathematical research problem. However, the Polymath model required human moderators to review and integrate the contributions of the participants, which clearly would not scale to the ETP which required the verification of over twenty million mathematical statements. Instead, the ETP was centered around a GitHub repository in which the formal mathematical contributions had to be entered in the proof assistant language *Lean*, where they could be automatically verified. In this respect, the ETP was more similar to the recently concluded Busy Beaver Challenge⁶, which was a similarly crowd-sourced project that computed the fifth Busy Beaver number $BB(5)$ to be 47 176 870 through an analysis of about 180 million Turing machines, with the halting analysis being verified in a variety of computer languages, with the final formal proof written in the proof assistant language *Coq* [43, 44]. One of the aims of the ETP was to explore potential workflows for such collaborative, formally verified mathematical research projects that could serve as a model for future projects of this nature.

Secondary aims of the ETP included the possibility of discovering unusually interesting equational laws, or new experimental observations about such laws, that had not previously been noticed in the literature; and to develop benchmarks to assess the performance of automated theorem provers and other AI tools.

⁵We improved such a proof to make it human-readable, see the blueprint of the ETP.

⁶<https://bbchallenge.org/>

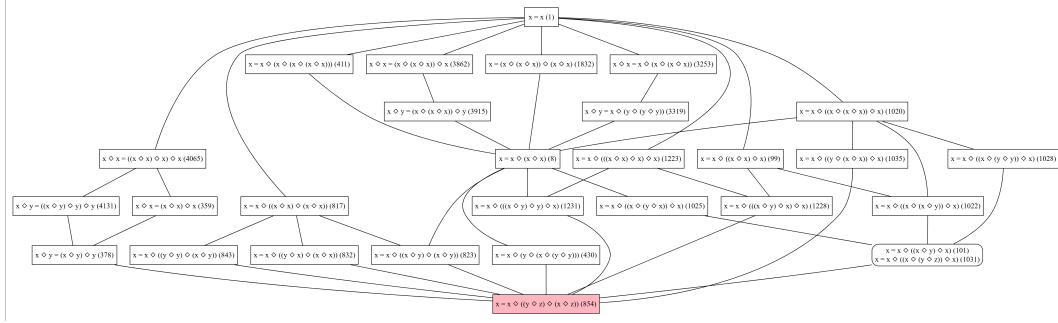


FIGURE 1. A Hasse diagram of all the equational laws implied by (E854) (for unrestricted magmas). An edge in this diagram indicates that the lower equation implies the higher one. Rounded rectangles indicate groups of equivalent laws. This graph was produced by the visualization tool *Graphiti*, which was developed for this project.

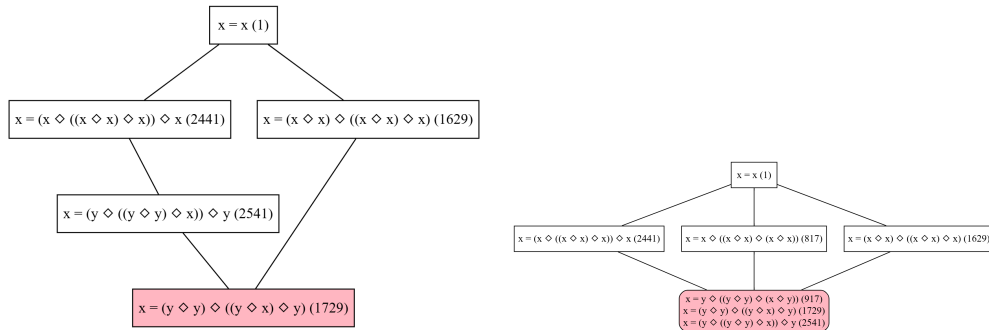


FIGURE 2. A Hasse diagram of all the equational laws implied by (E1729), both for unrestricted magmas (left) and finite magmas (right). Note the slightly larger number of implications in the latter.

1.3. Outcomes. The ETP achieved almost all of its primary objectives, with all of the 22 033 636 implications $E \vdash E'$ and non-implications $E \not\vdash E'$ magmas formalized in the proof assistant language *Lean*, and can be found on the ETP GitHub repository. See Figure 1, Figure 2 for some small fragments of the implication graphs produced. The 4694 laws organized into 1415 equivalence classes, with by far the largest class being the class of 1486 equations equivalent to the trivial law E2.

For the finite implication graph $E \vdash_{\text{fin}} E'$, we could similarly formalize all but two implications. Specifically, we were unable to obtain either a human-readable or formalized proof or disproof of the implication $E677 \vdash_{\text{fin}} E255$ (or its equivalent dual $E2910 \vdash_{\text{fin}} E47$), despite extensive efforts from the participants of the project; we tentatively conjecture this implication to be false (i.e., that there exists a finite magma obeying (E677) but not (E255)), but the refutation appears to be “immune” to most of the techniques that we developed for the project. (We were however able to establish that the corresponding implication $E677 \vdash E255$ for arbitrary magmas was false, using the greedy construction discussed in Section 5.5.)



FIGURE 3. The longest chain of implications (length 15) between inequivalent laws in the implication graph.

Of the 22 033 636 possible implications $E \vdash E'$, 8178279 (or 37.12%) would end up being true; for an additional set of either 820 or 822 pairs E, E' , the weaker implication $E \vdash_{\text{fin}} E'$ also held. To establish such positive implications $E \vdash E'$ or $E \vdash_{\text{fin}} E'$, the main techniques used were as follows:

- A very small number of positive implications were established and **formalized by hand**, mostly through direct rewriting of the laws; but this approach would not scale to the full project.
- **Simple rewriting rules**, for instance based on the observation that any law of the form $x \simeq f(y, z, \dots)$ was necessarily equivalent to the trivial law (E2), could already reduce the size of potential equivalence classes by a significant fraction. We discuss this method in Section 6.1.
- The preorder axioms for \vdash , as well as the “duality” symmetry of the preorder with respect to replacing a magma operation $x \diamond y$ with its opposite $x \diamond^{\text{op}} y := y \diamond x$, can be used to significantly cut down on the number of implications that need to be proven explicitly; ultimately, only 10657 (0.05%) of the positive implications needed a direct proof.
- To obtain additional implications for finite magmas, heavy reliance was made on the fact that for functions $f: M \rightarrow M$ on a finite set M , surjectivity was equivalent to injectivity. Some more sophisticated variants of this idea can lead to additional implications; see Section 5.1.

- **Automated Theorem Provers** (ATP) could be deployed at extremely fast speeds to establish a complete generating set of positive implications; see Section 7.

More challenging were the 13 855 357 (62.88%) implications that were false, $E \not\vdash E'$, and particularly the slightly smaller set of 13 854 535 or 13 854 537 implications that were false even for finite magmas, $E \not\vdash_{\text{fin}} E'$. Here, the range of techniques needed to refute such implications were quite varied, and may be of independent interest:

- **Syntactic methods**, such as observing a “matching invariant” of the law E that was not shared by the law E' , could be used to obtain some refutations. For instance, if both sides of E had the same order, but both sides of E' did not, this could be used to syntactically refute $E \vdash E'$. Similarly, if the law E was confluent, enjoyed a complete rewriting system, or otherwise permitted some understanding of the free magma associated to that law, one could decide the assertions $E \vdash E'$ for all possible laws E' , or at least a significant fraction of such laws. We discuss these methods, and the extent to which they can be automated in Section 6.
- **Small finite magmas**, which can be described explicitly by multiplication tables, could be tested by brute force computations to provide a large number of finite counterexamples to implications, or by ATP-assisted methods. See Section 5.1.
- **Linear models**, in which the magma operation took the form $x \diamond y = ax + by$ for some (commuting or non-commuting) coefficients a, b , allowed for another large class of counterexamples to implications, which could be automatically scanned for either by brute force or by Grobner basis type calculations; many of these examples could also be made finite. See Section 5.2.
- **Translation invariant models**, in which the magma operation took the form $x \diamond y = x + f(y - x)$ on an additive group, or $x \diamond y = xf(x^{-1}y)$ on a non-commutative group, reduce matters to analyzing certain functional equations; see Section 5.3.
- **Greedy methods**, in which either the multiplication table $(x, y) \mapsto x \diamond y$ or the function f determining a translation-invariant model are iteratively constructed by a greedy algorithm subject to a well-chosen ruleset, were effective in resolving many implications not easily disposed of by preceding methods. See Section 5.5.
- Starting with a simple base magma \mathcal{M} obeying both E and E' , and either **enlarging** it to a larger magma \mathcal{M}' containing \mathcal{M} as a submagma, **extending** it to a magma \mathcal{N} with a projection homomorphism $\pi : \mathcal{N} \rightarrow \mathcal{M}$, or *modifying the multiplication table* on a small number of values, also proved effective when combined with greedy methods or with a “**magma cohomology**” construction. See Section 5.6.
- To each equation E one can associate a “**twisting semigroup**” S_E . If S_E is larger than $S_{E'}$, then this can often be used to disprove the implication $E \vdash E'$; see Section 5.4.
- Some **ad hoc models** based on existing mathematical objects, such as infinite trees, rings of polynomials, or “Kisielewicz models” utilizing the prime factorization of the natural numbers, could also handle some otherwise difficult cases. In some cases, the magma law induced some relevant and familiar structures, such as a directed graph or a partial order, which also helped guide counterexample constructions. We will not detail these diverse examples here, but refer the reader to the ETP blueprint for more discussion.

- **Automated theorem provers** were helpful in identifying which simplifying axioms could be added to the magma without jeopardizing the ability to refute the desired implication $E \vdash E'$ or $E \vdash_{\text{fin}} E'$.

While the vast majority of negative implications could be quickly resolved by one of the above techniques, either with human input or in a completely automated fashion, there were perhaps two dozen such negative results that required quite delicate and *sui generis* constructions. The hardest such implication, E1729 $\not\vdash$ E817, took several months to establish and then formalize (using a combination of many of the above constructions), with the final proof in *Lean* requiring just over 4000 dedicated lines of code from multiple contributors.

In the course of completing the implication graph, some interesting new algebraic structures were discovered. One such example concerns the magmas obeying (E1485), which we refer to as *weak central groupoids* as they contain the central groupoids (obeying (E168)) as a subclass. In [19] it was observed that all finite central groupoids have order equal to a perfect square n^2 ; empirically, we have found that finite weak central groupoids always have order n^2 or $2n^2$, although we have no rigorous proof of this claim; they also have a graph-theoretic interpretation analogous to the interpretation of central groupoids as digraphs with the unique path property. For these and other observations we refer the reader to the blueprint of the ETP.

The objective of using the data from the ETP to establish well-calibrated benchmarks to evaluate ATPs remains an interesting open problem; the participants of this project did not have the required expertise to develop and test such benchmarks to the standards expected in the area. However, in Section 7 we present a more informal “field report” of our experiences using ATPs in the project, in the hope that this will provide some useful guidance to other researchers seeking to incorporate ATPs into their own research.

1.4. Further directions. While the primary objective of the ETP was being completed, some additional related results were generated as spinoffs. Specifically:

- In the blueprint on the ETP web site, we report some partial progress on we report on classifying which of the 57882 distinct laws of order 5 are equivalent to the singleton law (E2), either with or without the requirement that the magma be finite.
- In Section 9 we report on classifying the laws of order 8 that are equivalent to the Higman-Neumann law (E42323216).
- In Section 10 we report on preliminary experiments on using machine learning to determine to what extent the implication graph can be predicted by a neural network.

2. NOTATION AND MATHEMATICAL FOUNDATIONS

If $\mathcal{M} = (M, \diamond)$ is a magma, we define the left and right multiplication operators $L_a, R_a: M \rightarrow M$ for $a \in M$ by the formula

$$(1) \quad L_x y = R_y x := x \diamond y.$$

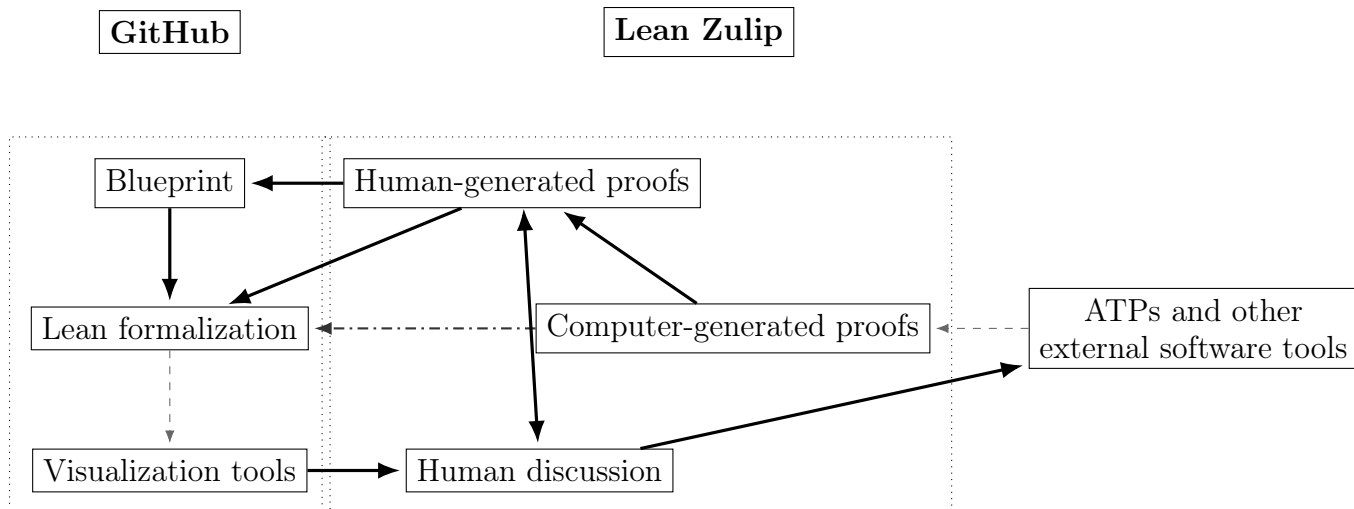


FIGURE 4. Some of the main dynamics in which proofs were generated, discussed within the Lean Zulip channel and then formalized in the Github repository. Boldface arrows indicate human activities, such as proposing an automated attack on outstanding implications, converting a computer-generated proof into a human-readable format, formalizing a human readable proof directly, or first creating a more precise blueprint for other collaborators to work on. Dashed arrows indicate fully automated processes, while the partly dashed line indicated a semi-automated process requiring human supervision. **figure out where to put this image and discuss it**

We also define the squaring operator $S: M \rightarrow M$ by

$$(2) \quad Sx := x \diamond x = L_x x = R_x x.$$

A *homomorphism* $f: \mathcal{M} \rightarrow \mathcal{M}'$ between two magmas $\mathcal{M} = (M, \diamond)$, $\mathcal{M}' = (M', \diamond')$ is a function $f: M \rightarrow M'$ such that $f(x \diamond y) = f(x) \diamond' f(y)$ for all $x, y \in M$. An *isomorphism* is a homomorphism that is invertible (which implies that the inverse is also a homomorphism). An *endomorphism* is a homomorphism from a magma to itself.

If X is an alphabet, we let \mathcal{M}_X denote the free magma generated by X , thus an element of \mathcal{M}_X is either a letter in X , or of the form $w_1 \diamond w_2$ with $w_1, w_2 \in \mathcal{M}_X$. Every function $f: X \rightarrow M$ into a magma $\mathcal{M} = (M, \diamond)$ extends to a unique homomorphism $\varphi_f: \mathcal{M}_X \rightarrow \mathcal{M}$. Formally, an equational law with some indeterminates in X can be written as $w_1 \simeq w_2$ for some $w_1, w_2 \in X$; a magma $\mathcal{M} = (M, \diamond)$ then obeys this law if and only if $\varphi_f(w_1) = \varphi_f(w_2)$ for all $f: X \rightarrow M$. We also define the order of a word $w \in \mathcal{M}_X$ to be the number of occurrences of \diamond in the word, thus letters in X are of order 0, and the order of $w_1 \diamond w_2$ is the sum of the orders of w_1, w_2 .

A *theory* is a collection Γ of equational laws; we say that a magma \mathcal{M} *satisfies* a theory, and write $\mathcal{M} \models \Gamma$, if every law in Γ is obeyed by \mathcal{M} . If E is an equational law, we write $\Gamma \vdash E$ if every magma that satisfies Γ also satisfies E . A *free magma* $\mathcal{M}_{X,\Gamma}$ for such a theory Γ and an alphabet X is a magma satisfying Γ together with a map $\iota_{X,\Gamma}: X \rightarrow \mathcal{M}_{X,\Gamma}$ which is universal in the sense that every function $f: X \rightarrow \mathcal{M}$ to a magma \mathcal{M} satisfying Γ uniquely

determines a homomorphism $\varphi_{f,\Gamma}: \mathcal{M}_{X,\Gamma} \rightarrow \mathcal{M}$ such that $\phi_{f,\Gamma} \circ \iota_{X,\Gamma} = f$. This magma is unique up to isomorphism; a canonical way to construct it is as the quotient $\mathcal{M}_X / \sim_\Gamma$ of the free magma \mathcal{M}_X by the equivalence relation \sim_Γ given by declaring $w \sim_\Gamma w'$ if $\Gamma \vdash w \simeq w'$ [3, Theorem 3.5.6]. If $\Gamma = \{E\}$ consists of a single law E , we write $\mathcal{M}_{X,E}$, \sim_E , $\varphi_{f,E}$ for $\mathcal{M}_{X,\{E\}}$, $\sim_{\{E\}}$, $\varphi_{f,\{E\}}$ respectively.

In general, the free magma $\mathcal{M}_{X,\Gamma}$ is difficult to describe in a tractable form, but for some theories, one has a simple description. We give two simple examples here:

Example 2.1 (Commutative and associative free magma). The free magma $\mathcal{M}_{X,\{E43,E4512\}}$ for the commutative law (E43) and the associative law (E4512) is the free abelian semigroup generated by X (with $\iota_{X,\{E43,E4512\}}$ the obvious embedding map).

Example 2.2 (Left-absorptive free magma). The free magma $\mathcal{M}_{X,\{E4\}}$ for the left-absorptive law (E4) is the magma with carrier X and operation $x \diamond y = x$ (with $\iota_{X,E4}$ the identity).

Every magma \mathcal{M} has an opposite \mathcal{M}^{op} , which has the same carrier but the opposite operation $x \diamond^{\text{op}} y := y \diamond x$. A magma \mathcal{M} obeys an equational law E if and only if its opposite \mathcal{M}^{op} obeys the dual law E^* , defined by reversing the all operations. For instance, the dual of $x \diamond y \simeq x \diamond (y \diamond z)$ (E327) is $y \diamond x \simeq (z \diamond y) \diamond x$, which in our numbering system we rewrite in normal form as $x \diamond y \simeq (z \diamond x) \diamond y$ (E395).

We then see that the implication graph has a duality symmetry: given two equational laws E_1, E_2 , we have $E_1 \vdash E_2$ if and only if $E_1^* \vdash E_2^*$.

3. FORMAL FOUNDATIONS

TODO: expand this sketch.

All proofs in the ETP were ultimately formalized in the proof assistant language *Lean*, though in many cases the proofs were first written either in an informal human document, which was then incorporated into the human-readable *blueprint*⁷ that accompanied the formalization. Many of the computer assisted proofs were also first generated as computer output from a source other than *Lean*, such as an ATP, and later converted to a *Lean* proof by a separate program custom-written for this task.

The project relied on *Lean*'s extensive *Mathlib* library, for instance to provide support for algebraic concepts such as the free group that arose in some of the more difficult constructions. The concept of a magma could be modeled by the *Mathlib* class `Mul`, however we chose early in the project to define a custom magma class `Magma` instead, as for some magma constructions the magma operation (which we denoted \diamond) was distinct from an existing multiplication structure $*$ on the same carrier. Most components of the *Lean* codebase were placed in namespaces to avoid collisions with each other, and with *Mathlib*.

Equational laws in the project were implemented both syntactically - as a structure `LawX` containing two words in a free group - as well as semantically, as a predicate `EquationX`

⁷<https://github.com/PatrickMassot/leanblueprint>

```
@[equational_result]
theorem _root_.Equation1437_not_implies_Equation4269 :  $\exists$  (G : Type) ( _ : Magma G),
Equation1437 G  $\wedge$   $\neg$  Equation4269 G := by
  use  $\mathbb{N} \times \text{Fin } 3$ , <op>
  constructor
  · intro x y z
    simp [op, add_assoc]
  · simp only [not_forall, op]
    use (0, 0), (2, 0)
  decide
```

FIGURE 5. A sample proof of a formalized implication, in this case that $E_{1437} \not\vdash E_{4269}$.

```
@[equational_result]
theorem «Facts from All4x4Tables [[1,2,3,4,5,0],[4,1,2,5,0,3],[3,0,5,2,1,4],
[0,5,4,3,2,1],[5,4,1,0,3,2],[2,3,0,1,4,5]]» :
 $\exists$  (G : Type) ( _ : Magma G) ( _ : Finite G), Facts G [1316, 2863] [411, 680, 817, 1020,
1426, 2035, 2441, 2644, 2853, 2855, 2865, 2872, 2947, 3050, 3253, 3456, 4270,
4283, 4290, 4380, 4598, 4605, 4656] :=
<Fin 6, «All4x4Tables [[1,2,3,4,5,0],[4,1,2,5,0,3],[3,0,5,2,1,4],[0,5,4,3,2,1],
[5,4,1,0,3,2],[2,3,0,1,4,5]]», Finite.of_fintype _, by decideFin!>
```

FIGURE 6. A computer generated `Facts` theorem, using an explicit finite magma of order 6 to refute several implications at once.

that could be applied to a magma. Here X is the number assigned to the law. The semantic formulation (`EquationX`) was more convenient for proving or refuting specific implications, while the syntactic formulation (`LawX`) was preferred for implementing metatheorems, such as the use of duality between laws. *Lean*'s metaprogramming features proved to be vital relate both representations. A custom command, `equation`, was created for specifying equational laws. Elaborating the `equation` command generated both `EquationX` and `LawX` definitions from this description, as well as theorems relating them to each other. A similar construction was used to generate dual laws, where the dual law was given explicitly for simplicity.

To facilitate the automatic generation of an implication graph from the *Lean* codebase, a custom `@[equational_result]` tag was formed to attach to propositions in *Lean* to indicate that they were proving or refuting one or more implications; see Figure 5. A `conjecture` keyword was also created for implications or refutations which we wished to identify as having an informal proof that had yet to be formalized in *Lean*.

A single construction of a magma could obey multiple laws E_1, E_2, \dots and not obey others E'_1, E'_2, \dots , leading to a large number of refutations of the form $E_i \not\vdash E'_j$. A custom `Facts` command was designed to organize such information efficiently; see Figure 6.

As an additional precaution against “exploit”-based proofs (such as those that might be contributed by an AI tool) *lean4checker* was used to ensure that no axioms were used in

Lean outside of a small trusted set. In particular, *Lean* tactics such as *native_decide* that relied on external tools were not permitted into the codebase.

This is my vague understanding of the situation - someone with experience should validate this.

Explicitly formalizing all 22 028 942 implications as theorems would lead to an infeasible compilation time in *Lean*. Instead, a reduced generating set of 10 657 positive implications and 586 925 negative implications were formalized, with the latter in turn mostly organized into a smaller number of **Facts** theorems as discussed above. The extension of these results to the rest of the implication graph via transitivity and duality is currently done by programs external to *Lean*, although in principle one could create an “end-to-end theorem” which completely establishes the implication graph within *Lean*.

Some lemmas generated in the project were suitable for upstreaming back to *Mathlib*, as well as several technical improvements to the *LeanBlueprint* software.

4. PROJECT MANAGEMENT

Shreyas Srinivas and Pietro Monticone have volunteered to take the lead on this section.

This project is, among other things, an experiment on how to organise large scale collaborations for mathematical work. In this section, we describe several aspects of the mechanics of the collaborative effort.

4.1. Problems of scale in mathematical collaboration. In order to understand the difficulties of scaling that arise in large scale collaborations, it helps to revisit how traditional mathematical collaborations work and understand why they may not scale. Every collaboration is unique, and we cannot imagine that any universal template exists. However, some general patterns can be observed in any mathematical collaboration. Usually, a small number of contributors, usually under ten, who may know each other, join forces to tackle some class of problems. Typically the collaborators are academics who share substantial amounts of common knowledge. They discuss the problem at hand together, typically with some shared written medium such as a whiteboard. After several rounds of discussion and refinement, different members of the collaboration come up with different pieces of the solution and discuss how they fit together into a whole. Along the way, each collaborator writes out their specific contributions and reviews that of others. After several iterations of this process they synthesize the results into a single paper. At this point, the collaborators are reasonably confident about the correctness of their work, including theorem statements and proofs, and submit the paper for peer review. Of course, there are many variations to this general template, but the basic cycle of discuss, solve, write, cross-check, and revise process is always present in some form or another.

Discuss how the above breaks down when a large and diverse collection of collaborators work over the internet. Consider citing the EMS webpage for code of practice for joint responsibility rules. Start the project organisation setup in the next subsection

Discuss topics such as:

- Project generation from template
- GitHub issue management with labels and task management dashboard
- Continuous integration (builds, blueprint compilation, task status transition)
- Pre-push git hooks
- Use of blueprint (small note, see #406: blueprint chapters should be given names for stable URLs)
- Use of Lean Zulip (e.g. use of polls)

Maybe give some usage statistics, e.g. drawing from https://github.com/teorth/equational_theories/actions/metrics/usage

Mention that FLT is also using a similar workflow.

4.2. Handling Scaling Issues. Also mention some early human-managed efforts ("outstanding tasks", manually generated Hasse diagram, etc.) which suffices for the first one or two days of the project but rapidly became unable to handle the scale of the project.

Mention that some forethought in setting up a GitHub organizational structure with explicit admin roles etc. may have had some advantages if we had done so in the planning stages of the project, but it was workable without this structure (the main issue is that a single person – Terry – had to be the one to perform various technical admin actions).

Use of transitive reduction etc. to keep the Lean codebase manageable. Note that the project is large enough that one cannot simply accept arbitrary amounts of Lean code into the codebase, as this could make compilation times explode. Also note somewhere that transitive completion can be viewed as directed graph completion on a doubled graph consisting of laws and their formal negations.

Technical debt issues, e.g., complications stemming from an early decision to make `Equations.lean` and `AllEquations.lean` the ground truth of equations for other analysis and visualization tools, leading to the need to refactor when `AllEquations.lean` had to be split up for performance reasons.

Note that the "blueprint" that is now standard for guiding proof formalization projects is a bit too slow to keep up with this sort of project that is oriented instead about proving new results. Often new results are stated and formalized without passing through the blueprint, which is then either updated after the fact, or not at all. So the blueprint is more of an optional auxiliary guiding tool than an essential component of the workflow.

4.3. Other Design Considerations. Explain what "trusting" Lean really means in a large project. Highlight the kind of human issues that can interfere with this and how use of tools like external checkers and PR reviews by people maintaining the projects still matters. Provide guidelines on good practices (such as branch protection or watching out for spurious modifications in PRs, for example to the CI). Highlight the importance of following a proper process for discussing and accepting new tasks, avoiding overlaps etc. These issues are less likely to arise in projects with one clearly defined decision maker as in this case, and more likely to arise when the decision making has to be delegated to many maintainers.

Note that despite the guarantees provided by Lean, non-Lean components still contained bugs. For instance, an off-by-one error in an ATP run created a large number of spurious conjectures, and some early implementations of duality reductions (external to Lean) were similarly buggy. "Unit tests", e.g., checking conjectured outputs against Lean-validated outputs, or by theoretical results such as duality symmetry, were helpful, and the Equation Explorer visualization tool also helped human collaborators detect bugs.

Meta: documenting thoughts for the future record is quite valuable. It would be extremely tedious to try to reconstruct real-time impressions long after the fact just from the GitHub commit history and Zulip chat archive.

4.4. Maintenance. Describe the role of maintainers and explain why they need to be conversant in the mathematics being formalised, as well as Lean. As such, the role of maintainers is often akin to postdocs or assistant profs in a research group who do some research of their own, but spend much of their time to guide others in performing their tasks, the key difference being that contributors are volunteers who choose their own tasks. Explain the tasks maintainers must perform. Examples:

- Reviewing proofs,
- Helping with proofs and theorem statements when people get stuck,
- Offering suggestions and guidance on how to produce shorter or more elegant proofs,
- Ensuring some basic standards are met in proof blocks that make proofs robust to upstream changes,
- Creating and maintaining CI processes,
- Responding to task requests,
- Evaluating theorem and definition formulations (for example unifying many theorem statements into one using FactsSyntax),
- Suggesting better ones where possible,
- Ensuring that there is no excessive and pointless overlap of content in different contributions

elaborate on what level of overlap was permissible and what we consider excessive

5. COUNTEREXAMPLE CONSTRUCTIONS

In this section we collect the various techniques developed in the ETP to construct counterexamples to various implications $E \vdash E'$.

5.1. Finite magmas. A finite magma \mathcal{M} of order n can be labeled by the carrier $\{1, \dots, n\}$ and described by specifying the multiplication table $\diamond: \{1, \dots, n\} \times \{1, \dots, n\} \rightarrow \{1, \dots, n\}$. By generating a list of all the equational laws E_j , $j = 1, \dots, 4694$ obeyed by this magma, one can create refutations: if $\mathcal{M} \models E_j$ and $\mathcal{M} \not\models E_k$, then clearly $E_j \not\vdash_{\text{fin}} E_k$ and hence also $E_j \not\vdash E_k$. (As mentioned previously, these statements were organized in Lean using the **Facts** statement.) It is feasible to brute force over all $\sum_{n=2}^4 n^{n^2} \approx 4.3 \times 10^9$ non-trivial magmas of order at most 4 to obtain many refutations of this type. By performing brute force over all magmas up to order 4, a total of 13,632,566 implications (61.9% of all implications, and 96.3% of the false ones) can be refuted with 524 distinct magmas. Of these implications, 13,345,053 were refuted with 3×3 magmas, with the remaining 415,293 requiring 4×4 magmas. Performing this search took 165 CPU-hours.

However, it is not feasible to exhaustively search over the $5^{5^2} \approx 3 \times 10^{17}$ magmas of order 5, even after quotienting out by isomorphism and symmetry (which roughly saves a factor of $5! \times 2 = 240$). Randomly sampling such magmas did not produce significant refutations, as random magmas of order 5 tended to obey few laws, and the set of laws covered were usually also exhibited by smaller magmas. A more fruitful approach was to randomly sample from magmas with additional properties that encouraged satisfiability of a greater set of laws. These included linear and quadratic magmas (discussed below), and cancellative magmas. On the other hand, some classes of magmas, such as commutative magmas, ended up producing a disappointingly small number of additional refutations.

For specific refutations, it was sometimes possible to locate a finite example with an ATP, particularly if one also imposed additional axioms (e.g., an idempotence axiom $x = x \diamond x$) that one suspected would be useful; see Section 7 for further discussion. For medium-sized magmas (of order $n = 5, 6, 7, 8$), this appeared to be a more efficient approach than brute force exhaustion of all such magmas.

Discuss this further, perhaps give an example, or refer to the ATP section. See also the discussion threads “Counterexamples by enumerating words in quotient magmas” and “Using SAT solvers for model generation” threads (sorry, LaTeX is choking on the URLs for some reason).

It is a result of Kisielewicz [17] that every law En with $n \leq 4694$ is either equivalent to the singleton law E2, or else has a non-trivial finite model; in other words, the implications $En \vdash E2$ and $En \vdash_{\text{fin}} E2$ are equivalent for $n \leq 4694$. In fact our brute force search revealed that in the latter case there is always a model of order $2 \leq n \leq 5$, with the lone exception of (E1286) (and its dual (E2301)), for which the smallest non-trivial finite model was of order 7, as presented in Theorem 5.2 below. In fact, most of the 4694 laws either only had trivial models, or had an order 2 model, as shown in Table 1.

| Order of smallest non-trivial model | Number of laws |
|-------------------------------------|----------------|
| Trivial only | 1496 |
| 2 | 3136 |
| 3 | 32 |
| 4 | 14 |
| 5 | 14 |
| 7 | 2 |

TABLE 1. Number of laws of order at most 4 whose smallest non-trivial model (if any) is of a given size.

Remark 5.1. The laws obeyed by a given finite magma M need not be finitely axiomatizable. The smallest example⁸ is the three-element magma $\{0, 1, 2\}$ with $1 \diamond 2 = 1$, $2 \diamond 1 = 2 \diamond 2 = 2$, and $x \diamond y = 0$ for all other x, y [35]. It was also shown in [36] that “almost all” magmas M (in a certain precise sense) are idempotential, which implies that their laws are finitely axiomatizable, and all other finite magmas obeying these laws are isomorphic to powers of M .

5.2. Linear models. As it turns out, a particularly fruitful source of counterexamples is the class of *linear magmas*, where the carrier M is a ring (which may be commutative or non-commutative, finite or infinite), and the operation \diamond is given by $x \diamond y = ax + by$ for some coefficients $a, b \in M$; one can also generalize this slightly to *affine magmas*, in which the operation is given by $x \diamond y = ax + by + c$, but for simplicity we shall focus on linear magmas here. It is easy to see that in a linear magma, any word $w(x_1, \dots, x_n)$ of n indeterminates also takes the linear form

$$w(x_1, \dots, x_n) = \sum_{i=1}^n P_{w,i}(a, b)x_i$$

for some (possibly non-commutative) polynomial $P_{w,i}$ in a, b with integer coefficients. Thus, a linear magma will obey an equational law $w_1 \simeq w_2$ if and only if the pair (a, b) lies in the (possibly non-commutative) variety

$$(3) \quad V_{w_1, w_2}(M) := \{(a, b) \in M \times M : P_{w_1, i}(a, b) = P_{w_2, i}(a, b) \text{ for all } i\} \subset M^2.$$

As such, a necessary condition for such a law $w_1 \simeq w_2$ to entail another law $w'_1 \simeq w'_2$ is that one has the inclusion

$$V_{w_1, w_2}(M) \subset V_{w'_1, w'_2}(M)$$

for all rings M . For commutative rings, this criterion can be checked by standard Grobner basis techniques; in the noncommutative case one can use methods such as the diamond lemma [5].

Example 5.2 (Commutative counterexample). For the law $x = y \diamond ((x \diamond y) \diamond x) \diamond y$ (E1286), the variety Equation (3) can be computed to be

$$\{(a, b) \in M \times M : 1 = a + ba^3 + bab, 0 = a + ba^2b + b^2\}$$

while the variety for the idempotent law (E3) is

$$\{(a, b) \in M : a + b = 1\}.$$

⁸We thank Stanley Burris for these references.

Thus, to show that (E1286) does not entail (E3), it suffices to locate elements a, b of a ring M for one has $1 = a + ba^3 + bab$, $0 = a + ba^2b + b^2$, and $a + b \neq 1$. Here one can take a commutative example, for instance when $M = \mathbb{Z}/p\mathbb{Z}$ and $(p, a, b) = (11, 1, 7)$ or $(p, a, b) = (7, 6, 2)$.

Example 5.3 (Noncommutative counterexample). For the law $x = y \diamond ((y \diamond (x \diamond z)) \diamond z)$ (E1117), the variety Equation (3) can be computed to be

$$\{(a, b) \in M \times M : 1 = baba, 0 = a + ba^2, 0 = bab^2 + b^2\}$$

while the variety for $x = (x \diamond ((x \diamond x) \diamond x)) \diamond x$ (E2441) is

$$\{(a, b) \in M \times M : a^2 + aba^2 + abab + ab^2 + b = 1\}.$$

Observe that if $ba = -1$, then (a, b) automatically lies in the first set, and lies in the second set if and only if $(ab+1)(b-1) = 0$. One can then show that (E1117) does not imply (E2441) by setting $a = L$, $b = -R$ where L, R are the left and right shift operators respectively on the ring of integer-valued sequences $\mathbb{Z}^{\mathbb{N}}$. With some *ad hoc* effort one can convert this example into a less linear, but simpler (and easier to formalize) example, namely the magma with carrier \mathbb{Z} and operation $x \diamond y = 2x - \lfloor y/2 \rfloor$.

Remark 5.4. As essentially observed in [1], if there is a commutative linear counterexample to an implication $E \vdash E'$, then by the Lefschetz principle this counterexample can be realized in a finite field \mathbb{F}_q for some prime power q (and by the Chebotarev density theorem one can in fact take q to be a prime, so that the carrier is of the form $\mathbb{Z}/p\mathbb{Z}$ for some prime p), so that one also has $E \vdash_{\text{fin}} E'$. As such, we have found that an effective way to refute implications by the commutative linear magma method is to simply perform a brute force search over linear magmas $x \diamond y = ax + by$ in $\mathbb{Z}/p\mathbb{Z}$ for various triples (p, a, b) .

Discuss performance of this method.

On the other hand, the refutations obtained by non-commutative linear constructions need not have a finite model. For instance, consider the refutation $E1117 \not\vdash E2441$ from Theorem 5.3. The law (E1117) can be rewritten as $L_y R_z L_y R_z x = x$. This implies that R_z is injective and L_y is surjective for all y, z . For finite magmas \mathcal{M} , this then implies that the L_y, R_z are in fact invertible, and hence we have also $R_z L_y R_z L_y x = x$, which implies (E2441) by setting $x = y = z$. Thus, the refutation $E1117 \not\vdash E2441$ is “immune” to finite counterexamples.

Remark 5.5. One can also consider nonlinear magma models, such as quadratic models $x \diamond y = ax^2 + bxy + cy^2 + dx + ey + f$ in a cyclic group $\mathbb{Z}/N\mathbb{Z}$. For small values of N , we have found such models somewhat useful in providing additional refutations of implications $E \vdash_{\text{fin}} E'$ beyond what can be achieved by the linear or affine models. However, as the polynomials associated to a word $w(x_1, \dots, x_n)$ tend to be of high degree (exponential in the order of the word), it becomes quite rare for such models to obey a given equation E when N is large.

Remark 5.6. One can also consider the seemingly more general linear model $x \diamond y = ax + by$, where the carrier M is now an abelian group, and a, b act on M by homomorphisms, that is to say that they are elements of the endomorphism ring $\text{End}(M)$. However, this leads to exactly the same varieties Equation (3) (where M is now replaced by the endomorphism ring

$\text{End}(M)$) and so does not increase the power of the linear model for the purposes of refuting implications.

Give some statistics of what proportion of refutations can be resolved by linear models.

On the other hand, there are certainly some refutations $E \not\vdash E'$ of implications that are “immune” to both commutative and non-commutative models, in the sense that all such models that obey E , also obey E' . One such example is the refutation $\text{E1485} \vdash \text{E151}$, which we discuss further in Section 5.4 below.

5.3. Translation-invariant models. It is natural to look for counterexamples amongst magmas that obey a large number of symmetries. One such class of counterexamples are *translation-invariant models*, in which the carrier M is a group, and the left translations of this group form isomorphisms of the magma M . In the case of an abelian group $M = (M, +)$, such models take the form

$$(4) \quad x \diamond y = x + f(y - x)$$

for some function $f: M \rightarrow M$; in the case of a non-abelian group $M = (M, \cdot)$, such models instead take the form

$$(5) \quad x \diamond y = xf(x^{-1}y).$$

For such models, the verification of an equational law in n variables corresponds to a functional equation for f in $n - 1$ variables, as the translation symmetry allows one to normalize one variable to be the identity (say). This can simplify an implication to the point where an explicit counterexample can be found. These functional equations are trivial to analyze when $n = 1$. For $n = 2$, they are not as trivial, but still quite tractable, and has led to several refutations in practice. The method does not appear to be particularly effective for $n > 2$ due to the complexity of the functional equations.

Example 5.7 (Abelian example). For the law $x \simeq (x \diamond y) \diamond ((x \diamond y) \diamond y)$ (E1648), we apply the abelian translation-invariant model Equation (4) with $y = x + h$ to obtain

$$\begin{aligned} x \diamond y &= x + f(h) \\ (x \diamond y) \diamond y &= x + f(h) + f(h - f(h)) \\ (x \diamond y) \diamond ((x \diamond y) \diamond y) &= x + f(h) + f(f(h - f(h))) \end{aligned}$$

so that this law obeys (E1648) if and only if the functional equation

$$f(h) + f(f(h - f(h))) = 0$$

holds for all $h \in M$. Similarly, the law $x \simeq (x \diamond (x \diamond y)) \diamond y$ (E206) is obeyed if and only if

$$f(f(h)) + f(h - f(f(h))) = 0$$

for all $h \in M$. One can now check that the function $f: \mathbb{Z} \rightarrow \mathbb{Z}$ defined by $f(h) := -\text{sgn}(h)$ (thus $f(h)$ equals -1 when h is positive, $+1$ when h is negative, and 0 when h is zero) obeys the first functional equation but not the second, thus establishing that $\text{E1648} \not\vdash \text{E206}$.

Example 5.8 (Non-abelian example). We now obtain the opposite refutation $E_{206} \not\vdash E_{1648}$ to Theorem 5.7 using the non-abelian translation-invariant model. By similar calculations to before, we now seek to find a function $f: M \rightarrow M$ on a non-abelian group (M, \cdot) that obeys the functional equation

$$(6) \quad f(f(h))f(f(f(h))^{-1}h) = 1$$

for all $h \in M$, but fails to obey the functional equation

$$(7) \quad f(h)f(f(f(h))^{-1}h) = 1$$

for at least one $h \in M$. Now take M to be the group generated by three generators a, b, c subject to the relations $a^2 = b^2 = c^2 = 1$, or equivalently the group of reduced words in a, b, c with no adjacent letters in the word equal. We define

$$f(1) = 1, f(a) = b, f(b) = c, f(c) = a$$

and then $f(aw) = a$ for any non-empty reduced word w not starting with a , and similarly for b and c . The equation (6) can be checked directly for $h = 1, a, b, c$. If $h = aw$ with w non-empty, reduced, and not starting with a , then $f(f(h))^{-1} = f(f(h)) = b$ and $f(f(f(h))^{-1}h) = f(baw) = b$, giving (6) in this case, and similarly for cyclic permutations. Meanwhile, (7) can be checked to fail for $h = a$.

Remark 5.9. The construction in Theorem 5.8 also has the following more “geometric” interpretation. The carrier M can be viewed as the infinite 3-regular tree, in which every vertex imposes a cyclic ordering on its 3 neighbors (for instance, if we embed M as a planar graph, we can use the clockwise ordering). For $x, y \in M$, we then define $x \diamond y$ to equal x if $x = y$. If y is instead a neighbor of x , we define $x \diamond y$ to be the next neighbor of x in the cyclic ordering. Finally, if y is distance two or more from x , we define $x \diamond y$ to be the neighbor of x that is closest to y . One can then check that this model obeys (6) but not (7).

Remark 5.10. These constructions are necessarily infinitary in nature, because (E206) and (E1648) can be shown to be equivalent for finite magmas. Indeed, (E206) can be written as $x = R_y L_x L_x y$, which implies that R_y is surjective, hence injective, on a finite magma; writing $x = R_y z$ we conclude that $R_y z = R_y L_{z \diamond y} L_{z \diamond y} y$ and hence $z = L_{z \diamond y} L_{z \diamond y} y$, giving (E1648). The opposite implication is similar (using (E1648) to show that R_y is injective, hence surjective), and is left to the reader.

Some refutations $E \not\vdash E'$ are “immune” by translation-invariant models, in that any translation-invariant model that obeys E , also obeys E' . One obstruction is that for such models, the squaring map S is necessarily an invertible map, since $Sx = x + f(0)$ in the abelian case and $Sx = xf(1)$ in the non-abelian case. On the other hand, adding the assumption of invertibility of squares can sometimes force the implication $E \vdash E'$ to hold. For instance, the commutative law $x \diamond (y \diamond y) \simeq (y \diamond y) \diamond x$ (E4482) for a square and an arbitrary element will imply the full commutative law (E43) for translation-invariant models due to the surjectivity of S , but does not imply it in general (as one can easily see by considering models where S is constant).

5.4. The twisting semigroup. Suppose one has a magma \mathcal{M} obeying a law E , that also enjoys some endomorphisms $T, U: \mathcal{M} \rightarrow \mathcal{M}$. Then one can “twist” the operation \diamond by T, U

to obtain a new magma operation

$$(8) \quad x \diamond' y := Tx \diamond Uy.$$

If one then tests whether this new operation \diamond' obeys the same law E as the original operation \diamond , one will find that this will be the case provided that T, U obey a certain set of relations. The semigroup generated by formal generators T, U with these relations will be called the *twisting semigroup* Twist_E of E . This can be best illustrated with some examples.

Example 5.11. We compute the twisting semigroup of $x \simeq (y \diamond x) \diamond (x \diamond (z \diamond y))$ (E1485). We test this law on the operation Equation (8), thus we consider whether

$$x = (y \diamond' x) \diamond' (x \diamond' (z \diamond' y))$$

holds for all $x, y, z \in M$. Substituting in Equation (8) and using the homomorphism property repeatedly, this reduces to

$$x = (T^2y \diamond T U x) \diamond (U T x \diamond (U^2 T z \diamond U^3 y)).$$

If we impose the conditions $TU = UT$, $T^2 = U^3$, then this equation would follow from (E1485) (with x, y, z replaced with TUx, T^2y, U^2Tz respectively). Thus the twisting semigroup $\text{Twist}_{\text{E1485}}$ of (E1485) is generated by two generators T, U subject to the relations $TU = UT = 1$, $T^2 = U^3$. This is a cyclic group of order 5, since the relations can be rewritten as $T^5 = 1$, $U = T^{-1}$.

Now consider $x \simeq (x \diamond x) \diamond (x \diamond x)$ (E151). Applying the same procedure, we arrive at

$$x = (T^2x \diamond T U x) \diamond (U T x \diamond U^2x)$$

so the twisting group $\text{Twist}_{\text{E151}}$ is generated by two generators T, U subject to the relations $TU = UT = T^2 = U^2 = 1$. This is a cyclic group of order 2, since the relations can be rewritten as $T^2 = 1$, $U = T$.

Suppose the twisting semigroup Twist_E is not a quotient of $\text{Twist}_{E'}$, in the sense that the relations that define $\text{Twist}_{E'}$ are not obeyed by the generators of Twist_E . Then one can often disprove the implication $E \vdash E'$ by attempting the following procedure.

- First, locate a non-trivial magma $\mathcal{M} = (M, \diamond)$ obeying the law E . Then the Cartesian power M^{Twist_E} of tuples $(x_W)_{W \in \text{Twist}_E}$, with the pointwise magma operation, will also obey E .
- Furthermore, this Cartesian power admits two endomorphisms T, U defined by

$$T(x_W)_{W \in \text{Twist}_E} = (x_{WT})_{W \in \text{Twist}_E}; U(x_W)_{W \in \text{Twist}_E} = (x_{WU})_{W \in \text{Twist}_E},$$

which obey the relations defining Twist_E .

- We now twist the magma operation \diamond on M^{Twist_E} by T, U to obtain a new magma operation \diamond' defined by Equation (8), that will still obey law E .
- Because T, U will not obey the relations defining $\text{Twist}_{E'}$, it is highly likely that this twisted operation will not obey E' , thus refuting the implication $E \vdash E'$. If M and the twisting semigroup were finite, this approach should also refute $E \vdash_{\text{fin}} E'$.

For instance, a non-trivial finite model for (E1485) is given by the finite field \mathbb{F}_2 of two elements with the NAND operation $x \diamond y := 1 - xy$. If we twist \mathbb{F}_2^5 by the left shift $T(x_i)_{i=1}^5 =$

$(x_{i+1})_{i=1}^5$ and right shift $U(x_i)_{i=1}^5 = (x_{i-1})_{i=1}^5$, where we extend the indices periodically modulo 5, then the resulting operation

$$(x_i)_{i=1}^5 \diamond' (y_i)_{i=1}^5 := (1 - x_{i+1}y_{i-1})_{i=1}^5$$

on \mathbb{F}_2^5 will still obey (E1485), but will not obey (E151), thus showing that $E1485 \not\vdash_{\text{fin}} E151$ and hence $E1485 \not\vdash E151$. This particular implication does not seem to be easily establishable by any of the other methods discussed in this paper.

Report on how large the twisting semigroups are in practice, and how many implications can be refuted by this method.

5.5. Greedy constructions. We have found *greedy extension methods*, or *greedy methods* for short, to be a powerful way to refute implications, especially when the carrier M is allowed to be infinite. Such constructions have a long history in model theory, with possibly the earliest⁹ such construction due to Skolem [41]. A basic implementation of this method is as follows. To build a magma operation $\diamond: M \times M \rightarrow M$ that obeys one law E but not another E' , one can first consider *partial magma operations* $\diamond: \Omega \rightarrow M$, defined on some subset Ω of $M \times M$. Thus $x \diamond y$ is defined if and only if $(x, y) \in \Omega$. A magma operation is then simply a partial operation which is *total* in the sense that $\Omega = M \times M$. We say that a partial magma operation is *finitely supported* if Ω is finite.

In the language of first-order logic, a partial magma operation can also be viewed as a ternary relation $R(x, y, z)$ on M with the axiom that $R(x, y, z) \wedge R(x, y, z') \implies z = z'$ for all $x, y, z \in M$. The support Ω is then the set of (x, y) for which $R(x, y, z)$ holds for some (necessarily unique) z , which one can then take to be the definition of $z = x \diamond y$.

We say that one partial operation $\diamond': \Omega' \rightarrow M$ *extends* another $\diamond: \Omega \rightarrow M$ if Ω' contains Ω , and $x \diamond y = x \diamond' y$ whenever $x \diamond y$ (and hence $x \diamond' y$) are defined. Given a sequence $\diamond_n: \Omega_n \rightarrow M$ of partial operations, each of which is an extension of the previous, we can define the *direct limit* $\diamond_\infty: \bigcup_n \Omega_n \rightarrow M$ to be the partial operation defined by $x \diamond_\infty y := x \diamond_n y$ whenever $(x, y) \in \Omega_n$.

Abstractly, the greedy algorithm strategy can now be described as follows.

Theorem 5.12 (Abstract greedy algorithm). *Let E, E' be equational laws, and let Γ be a theory of first-order sentences regarding a partial magma operations $\diamond: \Omega \rightarrow M$ on a carrier M . Assume the following axioms:*

- (i) (*Seed*) *There exists a finitely supported partial magma operation $\diamond_0: \Omega_0 \rightarrow M$ satisfying Γ that contradicts E' , in the sense that there is some assignment of variables in E' in M such that both sides of E' are defined using \diamond_0 , but not equal to each other.*
- (ii) (*Soundness*) *If $\diamond_n: \Omega_n \rightarrow M$ is a sequence of partial magma operations obeying Γ with each \diamond_{n+1} an extension of \diamond_n , and the direct limit \diamond_∞ is total, then this limit obeys E .*

⁹We thank Stanley Burris for this reference.

(iii) (*Greedy extension*) If $\diamond: \Omega \rightarrow M$ is a finitely supported partial magma operation obeying Γ , and $a, b \in M$, then there exists a finitely supported extension $\diamond': \Omega' \rightarrow M'$ of \diamond to a possibly larger carrier M' , and also obeying Γ , such that $a \diamond' b$ is defined.

Then $E \not\vdash E'$.

We remark that this greedy method seems to be inherently infinitary in nature, and does not seem well adapted to refute finite magma implications $E \vdash_{\text{fin}} E'$.

Proof. We work on the countably infinite carrier \mathbb{N} . By embedding the finitely supported operation \diamond_0 from axiom (i) into \mathbb{N} , we can assume without loss of generality that \diamond_0 has carrier \mathbb{N} . By similar relabeling, we can assume in (iii) that $M' = M$ when $M = \mathbb{N}$, since any elements of $M' \setminus \mathbb{N}$ that appear in Ω' can simply be reassigned to natural numbers that did not previously appear in Ω . We well-order the pairs in $\mathbb{N} \times \mathbb{N}$ by (a_n, b_n) for $n = 1, 2, \dots$. Iterating (iii) starting from \diamond_0 , we can thus create a sequence of finitely supported magma operations $\diamond_0, \diamond_1, \dots$ on \mathbb{N} obeying Γ , with each \diamond_{n+1} an extension of \diamond_n , and $a_n \diamond_n b_n$ defined for all $n \geq 1$. Then the direct limit \diamond_∞ of these operations is total, and does not obey E' thanks to axiom (i). On the other hand, by axiom (ii) it obeys E , and the claim follows. \square

We refer to Γ as the *rule set* for the greedy extension method. To apply Theorem 5.12 to obtain a refutation $E \vdash E'$, we have found the following trial-and-error method to work well in practice:

1. Start with a minimal rule set Γ that has just enough axioms to imply the soundness property for the given hypothesis E .
2. Attempt to establish the greedy extension property for this rule set by setting $a \diamond' b$ equal to a new element $c \notin M$, and then defining additional values of \diamond' as necessary to recover the axioms of Γ' .
3. If this can be done in all cases, then locate a seed \diamond_0 refuting the given target E' , and STOP.
4. If there is an obstruction (often due to a ‘‘collision’’ in which a given operation $x \diamond' y$ is required to equal two different values), add one or more rules to Γ to avoid this obstruction, and return to Step 2.

As an example, we present

Proposition 5.13 (E73 does not imply E4380). *The law $x \simeq y \diamond (y \diamond (x \diamond y))$ (E73) does not imply $x \diamond (x \diamond x) \simeq (x \diamond x) \diamond x$ (E4380).*

Proof. To build a rule set Γ that will imply (E73) when total, a natural first choice would be the single rule

1. If $y \diamond (x \diamond y)$ is defined, then $y \diamond (y \diamond (x \diamond y))$ is defined and equal to x .

However, the greedy algorithm will fail just with this rule: if the partial operation has $x \diamond y$ and $z \diamond y$ both equal to some w for some $x \neq z$, then any attempt to assign a value to

$y \diamond (y \diamond w)$ will lead to a contradiction, as the above rule will force $y \diamond w$ to equal both x and z . Indeed, it is clear that (E73) forces all the right translation operators R_y to be injective. We therefore add this as an additional rule:

2. If $x \diamond y$ and $z \diamond y$ are defined and equal, then $x = z$.

To avoid some unwanted edge cases, it is also convenient to impose the additional rule

3. If $x \diamond y$ is defined, it is not equal to y .

Unlike Rule 2, this rule is not forced by (E73), but can be enforced as part of the greedy construction.

The ruleset clearly obeys the soundness axiom (ii) of Theorem 5.12. We now verify the greedy extension axiom (iii). Let Ω, a, b be as in that axiom. We may assume that $a \diamond b$ is undefined, since we are done otherwise. We adjoin a new element c to M to create M' , and set $a \diamond' b = c$. If we also have $b = d \diamond a$ for some d (unique by Rule 2, and only possible for $a \neq b$ by Rule 3), set $a \diamond' c = d$ (this is necessary to retain Rule 1). Of course, we also set $x \diamond' y = x \diamond y$ whenever $x \diamond y$ is already defined.

Since $c \notin M$, it is clear that \diamond' is a finitely supported partial magma operation on M' . It is also clear that \diamond' obeys Rule 2 and Rule 3. Now we case check Rule 1:

- Case 1: $x = c$ or $y = c$. Not possible since no left multiplication with c is defined.
- Case 2: $x \diamond' y = c$. Only possible when $x = a, y = b$, but then $y \diamond' (x \diamond' y)$ is undefined since $y = b \neq a$ if d is defined.
- Case 3: $y \diamond' (x \diamond' y) = c$. Only possible when $y = a$ and $x = d$, and holds in this case.
- Case 4: $x, y, x \diamond' y, y \diamond' (x \diamond' y) \neq c$: this case is covered by Rule 3 for \diamond .

To conclude, we need to locate a seed \diamond_0 obeying Rules 1,2,3 but contains a counterexample to (E4380). One simple example is the carrier $\{0, 1, 2, 3\}$ with $0 \diamond_0 0 = 1, 0 \diamond_0 1 = 2, 0 \diamond_0 2 = 0, 1 \diamond_0 0 = 3$. \square

This method is not guaranteed to halt in finite time, as there may be increasingly lengthy sets of rules one has to add to Γ to avoid collisions. However, in practice we have found many of the refutations that could not be resolved by simpler techniques to be amenable to this method (or variants thereof, as discussed below).

One can automate the above procedure by using ATPs (or SAT solvers) to locate new rules that are necessary and sufficient resolve any potential collision (and which, *a posteriori*, can be seen to be necessarily consequences of the law E). The seed-finding step (Step 3) is particularly easy to automate, and can also often be done by hand.

Describe performance of this automated method. Discuss the issue that some implications required a large SAT solver calculation that was difficult to formalize efficiently in Lean, prompting human-generated simplified proofs using smaller rulesets.

However, in some cases we have found it necessary to add “inspired” choices of rules that were not forced by the initial hypothesis E , but which simplified the analysis by removing problematic classes of collisions from consideration. We were unable to fully automate the process of guessing such choices; however, we found ATPs very useful for testing any proposed such guess. In particular, if an ATP was able to show that the existing ruleset, together with a proposed new rule A , implied E' , then this clearly indicated that one should not add A to the rule set Γ . Conversely, if an ATP failed to establish such an implication, this was evidence that this was a “safe” rule to impose.

We also found that human verification of the greedy extension property was a highly error-prone process, as the case analysis often included many delicate edge cases that were easy to overlook. Both ATPs and the Lean formalization therefore played a crucial role in verifying the human-written greedy arguments, often revealing important gaps in those arguments that required either minor or major revisions to the rule set.

The greedy method can also be combined with the translation-invariant method, both in abelian and non-abelian settings. For instance, we can modify the proof of Theorem 5.12 to obtain the following variant:

Theorem 5.14 (Non-commutative translation-invariant greedy algorithm). *Let F, F' be functional equations on groups, and let Γ be a theory of first-order sentences regarding a partial function $f: \Omega \rightarrow G$ on a group $G = (G, \cdot)$. Assume the following axioms:*

- (i) (Seed) *There exists a finitely supported partial function $f_0: \Omega_0 \rightarrow G$ satisfying Γ that contradicts F' , in the sense that there is some assignment of variables in F' in G such that both sides of F' are defined using f_0 , but not equal to each other.*
- (ii) (Soundness) *If $f_n: \Omega_n \rightarrow G$ is a sequence of partial functions obeying Γ with each f_{n+1} an extension of f_n , and the direct limit f_∞ is total, then this limit obeys F .*
- (iii) (Greedy extension) *If $f: \Omega \rightarrow G$ is a finitely supported partial function obeying Γ , and $a, b \in G$, then there exists a finitely supported extension $f': \Omega' \rightarrow G'$ of f to a possibly larger group G' , and also obeying Γ , such that $a \diamond' b$ is defined.*

Then $F \not\vdash F'$.

One can of course also develop an abelian analogue of the above theorem, in which $G = (G, +)$ and $G' = (G', +)$ are now required to be abelian. We can then give an alternate proof of Theorem 5.13 as follows:

Second proof of Theorem 5.13. (Sketch) The functional equations associated to (E73) and (E4380) are $f^2(h^{-1}f(h)) = h^{-1}$ and $f^2(1) = f(1)f(f(1)^{-1})$ respectively. We apply Theorem 5.14 with the following ruleset:

1. If $f(h^{-1}f(h))$ is defined, then $f^2(h^{-1}f(h))$ is defined and equal to h^{-1} .
2. If $h^{-1}f(h)$ and $k^{-1}f(k)$ are defined and equal, then $h = k$.
3. If $f(h)$ is defined, it is not equal to h .

Axiom (ii) is clear. To verify axiom (iii), we can assume $f(h)$ is undefined, then adjoin an element c freely to G to create a larger group G' , and set $f'(h) = c$. If $h = k^{-1}f(k)$ for some k (which is unique by Rule 2, and only possible for $h \neq 1$ by Rule 3), then also set $f'(c) = k^{-1}$. One can then check that axiom (iii) is obeyed. For axiom (i), take G to be a free cyclic group with one generator a , and set $f(1) = a$, $f(a) = a^3$, $f(a^3) = 1$, $f(a^{-1}) = a^3$ (say). \square

More complex (and *ad hoc*) variants of the greedy algorithm are possible. In some cases, we were not able to preserve the finitely supported nature of the partial operation or partial function, and needed to extend that partial object at an infinite number of values at each step. In other cases, one also had to add additional temporary data during the greedy process to record tasks that one wished to attend to at a later stage of the process, but could not handle immediately because it was awaiting some other operation to become well-defined. We will not attempt to survey all possible variants of this method here, but refer the reader to the ETP blueprint for further examples.

5.6. Modifying base models. A general technique that we have found useful in obtaining a refutation such as $E \not\vdash E'$ is to start with a simple base model $\mathcal{M} = (M, \diamond)$ that obeys both E and E' , and modify it in various ways to preserve E , but create a violation of E' . There are many such possible modifications, but three general ways that have proven effective are as follows:

- (i) Modify the magma operation $\diamond: M \times M \rightarrow M$ on a small set in order to violate E' , and then make further modifications as needed to recover E .
- (ii) Construct an *extension* \mathcal{N} of \mathcal{M} , equipped with a surjective magma homomorphism $\pi: \mathcal{N} \rightarrow \mathcal{M}$, and defined in terms of some additional data. Then solve for that data in such a way that \mathcal{N} obeys E but not E' .
- (iii) Construct an *enlargement* $\mathcal{M}' = (M', \diamond')$ of $\mathcal{M} = (M, \diamond)$, which is a magma that contains \mathcal{M} as a submagma. One needs to construct the multiplication table \diamond on $(M' \times M') \setminus (M \times M)$ in order to retain E but disprove E' .

One appealing case of (ii) that our project discovered, involving a “magma cohomology” analogous to (abelian) group cohomology, is that of an *affine* extension of a magma $\mathcal{G} = (G, \diamond_G)$ by another magma (M, \diamond_M) which is an abelian group M with a linear magma operation $s \diamond_M t := as + bt$ for some endomorphisms $a, b \in \text{End}(M)$. One can then consider extensions with carrier $G \times M$ and magma operation

$$(9) \quad (x, s) \diamond (y, t) := (x \diamond_G y, s \diamond_M t + f(x, y))$$

for some function $f: G \times G \rightarrow M$. If (M, \diamond_M) and (G, \diamond_G) already obey a law E , then this extension will also obey E if and only if f obeys a certain “cocycle equation”, which is a linear equation on f . One can then sometimes use linear algebra to locate an f that obeys the cocycle equation for one law E but not another E' , thus refuting the implication $E \vdash E'$. An example is as follows:

Proposition 5.15 (E1110 does not imply E1629). *The law $x \simeq y \diamond ((y \diamond (x \diamond x)) \diamond y)$ (E1110) does not imply $x \simeq (x \diamond x) \diamond ((x \diamond x) \diamond x)$ (E1629), even for finite magmas.*

Proof. (Sketch) Using the linear ansatz, we find that (E1110) has a model \mathcal{M} with carrier \mathbb{F}_5 (the finite field $\mathbb{Z}/5\mathbb{Z}$) with operation $x \diamond y = 3x - y$. We then apply the ansatz (9) with $G = M$. One then finds that this operation obeys (E1110) if $f: \mathbb{F}_5 \times \mathbb{F}_5 \rightarrow \mathbb{F}_5$ obeys the cocycle equation

$$3f(x, x) - 3f(y, 2x) - f(3y - 2x, y) + f(y, 3y - x) = 0$$

for all $x, y \in \mathbb{F}_5$, and obeys (E1629) if f obeys the cocycle equation

$$f(2x, 0) - f(2x, 2x) = 0$$

for all $x \in \mathbb{F}_5$. A routine symbolic algebra package computation reveals that the space of f that obeys the former equation is a six-dimensional vector space over \mathbb{F}_5 , which is not contained in the solution space of the latter equation, giving the claim. In fact, since these equations preserve the space of homogeneous polynomials of a fixed degree, one can use linear algebra to locate an example that is a homogeneous polynomial; one explicit choice is

$$f(x, y) = y^5 + xy^4 + x^2y^3 + 3x^3y^2 + 3x^4y_1$$

□

It may be of interest to develop this theory of “magma cohomology” further, for instance by defining higher order magma cohomology groups.

Now we give an example of how method (ii) can be combined with method (i).

Proposition 5.16 (E1659 does not imply E4315). $x \simeq (x \diamond y) \diamond ((y \diamond y) \diamond z)$ (E1659) *does not imply* $x \diamond (y \diamond x) \simeq x \diamond (y \diamond z)$ (E4315).

Proof. There are two simple models for (E1659): the model G with carrier $\mathbb{Z}/2\mathbb{Z}$ and operation $x \diamond y = x + 1$, and the model \mathcal{M} with carrier \mathbb{Z} and operation $x \diamond y = x$. Using the ansatz (9), one can soon discover that one obtains a magma operation $\diamond: (G \times M) \times (G \times M) \rightarrow G \times M$ with $f(0, 0) = f(1, 0) = 0$, $f(0, 1) = -1$, and $f(1, 1) = 1$. This model still obeys (E4315). However, we can create a modification \diamond' of \diamond as follows. We will seek to violate (E4315) at $x = (0, 0)$, $y = (0, 0)$, $z = (1, 0)$, thus we want

$$(0, 0) \diamond' ((0, 0) \diamond' (0, 0)) \neq (0, 0) \diamond' ((0, 0) \diamond' (1, 0)).$$

We have $(0, 0) \diamond (0, 0) = (1, 0)$ and $(0, 0) \diamond (1, 0) = (1, -1)$. One can try to force the counterexample by setting $(0, 0) \diamond' (1, 0)$ to equal $(0, 0)$ instead of $(1, -1)$. However, if this is the only change we make, then we no longer obey (E1659), since

$$(1, 0) \neq ((0, 0) \diamond' (1, 0)) \diamond' (((1, 0) \diamond' (1, 0)) \diamond (1, t))$$

for any $t \in \mathbb{Z} \setminus \{0\}$. But these are the only counterexamples created; and if one then sets $(0, 0) \diamond' (1, t) = (0, 0)$ for *all* $t \in \mathbb{Z}$, then one can check that the modified operation \diamond' now obeys (E1659) but not (E4315) as required. □

The specific law $x \simeq x \diamond ((y \diamond z) \diamond (x \diamond z))$ (E854) turned out to be somewhat “mutable”, in that one can often change a small number of entries in the multiplication table of a finite magma obeying this law, and also add an additional row and column to the table, in ways that preserve the law (E854). This renders this law suitable for using methods (i), (iii) to construct new models of this equation that refute various implications $E854 \vdash_{\text{fin}} E$, for

instance by starting with a model that already refuted some stronger law E' , and attempt to modify it (possibly with ATP assistance) by some combination of methods (i), (iii) to then violate E .

Another way to utilize (iii), which proved useful for laws that involved the squaring operator S , was to adopt a “squares first” approach in which one selected a base model $SM = (SM, \diamond)$ to serve as the set of squares, then extend it to a larger model \mathcal{M} with carrier $M = SM \uplus N$ by first determining what the multiplication map should be on the diagonal $\{(x, x) : x \in N\}$ (i.e., to determine the squaring map $S : N \rightarrow SM$), together with the values on the blocks $SM \times N$, $N \times SM$, and then finally resolve the remaining values on the $N \times N$ block. Often, versions of the greedy algorithm are useful for each of these stages of the construction. The precise details are quite technical, particularly for the law $x \simeq (y \diamond y) \diamond ((y \diamond x) \diamond y)$ (E1729), which was the last of the equations whose implications were settled by the ETP. We refer the reader to the ETP blueprint for further details.

6. SYNTACTIC ARGUMENTS

Many proofs or refutations of implications (or equivalences) between two equational laws E, E' can be obtained from the syntactic form of the equation. We discuss some techniques here that were useful in the ETP.

6.1. Simple rewrites. Many equational laws E' can be formally deduced from a given law E by applying the *Lean* `rw` tactic to rewrite E' repeatedly by some forward or backward application of E applied to arguments that match some portion of E . For instance, the commutative law (E43) clearly implies $x \diamond (y \diamond z) \simeq (y \diamond z) \diamond x$ (E4531) by a single such rewrite. A brute force application of such rewrite methods is already able to directly generate about 15 000 such implications, including many equivalences to the singleton law (E2) and the constant law (E46). After applying transitive closure, this generates about four million further such implications.

A simple observation that already generates a reasonable number of equivalences is that any equation of the form $x \simeq f(y, z, \dots)$ necessarily is equivalent to the trivial law $x \simeq y$, by transitivity; similarly, an equation of the form $f(x, y) \simeq g(z, w, \dots)$ implies $f(x, y) \simeq f(x', y')$; and so forth. Equivalences of this form were useful early in the project by cutting down the number of distinct equivalence classes of laws that needed to be studied.

6.2. Matching invariants. Fix an alphabet X . A *matching invariant* is an assignment $I : \mathcal{M}_X \rightarrow \mathcal{I}$ of an object $I(w) \in \mathcal{I}$ in some space \mathcal{I} to each word $w \in \mathcal{M}_X$ with the property that if an equational law $w_1 \simeq w_2$ has matching invariants $I(w_1) = I(w_2)$, then the same matching $I(w'_1) = I(w'_2)$ holds for any consequence $w'_1 \simeq w'_2$. In particular, if one law $I(w_1) = I(w_2)$ and $I(w'_1) \neq I(w'_2)$, then the law $w_1 \simeq w_2$ does not imply the law $w'_1 \simeq w'_2$.

A simple example of a matching invariant is the multiplicity $(n_x)_{x \in X}$ of variables of a word: if w_1, w_2 have all variables x appear the same number of times n_x in both words, then any rewriting of a word w using the law $w_1 \simeq w_2$ will preserve this property. Hence, if

w'_1, w'_2 do not have that each variable appear the same number of times in both words, then $w_1 \simeq w_2$ cannot imply $w'_1 \simeq w'_2$. For instance, the commutative law (E43) cannot imply the left-absorptive law (E4).

One source of matching invariants comes from the free magma $\mathcal{M}_{X,\Gamma}$ of a theory:

Proposition 6.1 (Free magmas and matching invariants). *Let Γ be a theory, and let $\iota_{X,\Gamma}: X \rightarrow \mathcal{M}_{X,\Gamma}$ be the map associated to the free magma $\mathcal{M}_{X,\Gamma}$ for that theory. Then the map $I: \mathcal{M}_X \rightarrow \mathcal{M}_{X,\Gamma}$ defined by $I(w) := \varphi_{\iota_{X,\Gamma}}(w)$ is an invariant.*

Proof. Suppose that $w_1 \simeq w_2$ entails $w'_1 \simeq w'_2$, and that $I(w_1) = I(w_2)$. For any $f: X \rightarrow \mathcal{M}_{X,\Gamma}$, the two maps $\varphi_f, \varphi_{f,\Gamma} \circ \varphi_{\iota_{X,\Gamma}}: \mathcal{M}_X \rightarrow \mathcal{M}_{X,\Gamma}$ are both homomorphisms that extend f , hence agree by the universal property of \mathcal{M}_X , as displayed by the following commutative diagram:

$$\begin{array}{ccccc}
 & & X & & \\
 & \swarrow & \downarrow \iota_{X,\Gamma} & \searrow f & \\
 \mathcal{M}_X & \xrightarrow{I = \varphi_{\iota_{X,\Gamma}}} & \mathcal{M}_{X,\Gamma} & \xrightarrow{\varphi_{f,\Gamma}} & \mathcal{M}_{X,\Gamma} \\
 & \searrow & \downarrow \varphi_f & \swarrow & \\
 & & & &
 \end{array}$$

In particular, the hypothesis $I(w_1) = I(w_2)$ implies that $\varphi_f(w_1) = \varphi_f(w_2)$ for all $f: X \rightarrow \mathcal{M}_{X,\Gamma}$; that is to say, the magma $\mathcal{M}_{X,\Gamma}$ obeys the law $w_1 \simeq w_2$, and hence also $w'_1 \simeq w'_2$ by hypothesis. In particular, $\varphi_{\iota_{X,\Gamma}}(w'_1) = \varphi_{\iota_{X,\Gamma}}(w'_2)$, which gives $I(w'_1) = I(w'_2)$ as required. \square

Example 6.2. If we take $\Gamma = \{\text{E4}\}$ to be the theory of the left-absorptive law (E4) as described in Theorem 2.2, then the matching invariant $I(w)$ produced by Theorem 6.1 is the left-most letter of the alphabet X appearing in the word; for instance $I((x \diamond y) \diamond z) = x$. Thus, for example, the left-absorptive law (E4) cannot imply the right-absorptive law (E5).

Example 6.3. If we take $\Gamma = \{\text{E43}, \text{Eq4512}\}$ to be the theory of the commutative law (E43) and the associative law (E4512), then by Theorem 2.1, the associated invariant $I(w) = \sum_{x \in X} n_x e_x$ is the formal sum of all the generators e_x appearing in the word w , in the free abelian semigroup generated by those generators. This recovers the preceding observation that the multiplicities $(n_x)_{x \in X}$ form a matching invariant.

Example 6.4. Let $n \geq 1$ be a positive integer, and consider the theory $\Gamma = \{\text{E43}, \text{E4512}, E_n\}$ consisting of the previous theory $\{\text{E43}, \text{E4512}\}$ together with the order- n law $L_x^y x = y$. One can check that the free magma $\mathcal{M}_{X,\Gamma}$ can be described as the free group of exponent n with generators $e_x, x \in X$, with associated map $\iota_{X,\Gamma}: x \mapsto e_x$. The associated matching invariant $I(w) = \sum_{x \in X} n_x e_x$ is essentially the multiplicities $(n_x \bmod n)_{x \in X}$ modulo n , which gives a slightly stronger criterion than the preceding matching invariant for refuting implications. For example, the cubic idempotent law $x \simeq (x \diamond x) \diamond x$ (E23) has matching invariants $e_x = 3e_x$ in the $n = 2$ case, and hence does not imply the idempotent law $x \simeq x \diamond x$ (E3) since $e_x \neq 2e_x$ in the $n = 2$ case.

In practice, we found that these invariants could be used to establish a significant fraction of the non-implications in the implication graph, although in most cases these non-implications

could also be established by other means, for instance through consideration of small finite counterexamples.

Remark 6.5. One can also obtain matching invariants from the free objects associated to theories that involve additional operations beyond the magma operation \diamond , such as an identity element or an inverse operation. We leave the precise generalization of Theorem 6.1 to such theories to the interested reader.

6.3. Confluence. We briefly recall the basics of rewrite theory necessary to our exposition, following mostly Baader and Nipkow [3], and generally omitting proofs when they can be found there.

We first work in the abstract taking an arbitrary set A , with a given equivalence relation over it which we denote \approx . We consider a relation R over A .

Definition 6.6. We write $a \rightarrow b$ if $a R b$ holds in A , and say that a *rewrites to* (or *reduces to*) b . We further define

- \rightarrow^+ as the transitive closure of R .
- \rightarrow^* as the reflexive transitive closure of R .
- \leftrightarrow^* as the reflexive transitive and symmetric closure of R .

We sometimes write $b \leftarrow a$, (resp. $b^* \leftarrow a$ etc) to mean $a \rightarrow b$ (resp. $a \rightarrow^* b$ etc), and chain notations, e.g. $b_1 \leftarrow a \rightarrow b_2$.

Example 6.7. Consider the relation over terms with a binary operation \diamond defined by:

$$\begin{aligned} x \diamond (y \diamond x) &\rightarrow y \\ (x \diamond y) \diamond x &\rightarrow y \end{aligned}$$

and the closure of this relation under substitution of the indeterminates x and y , as well as the *congruence closure* of the rules, that is $t \rightarrow t'$ if any subterm of t is related to the corresponding subterm of t' by the rules above.

Note that \leftrightarrow^* is an equivalence relation and the hope is for it to be equal to \approx , in order to deduce properties of the latter.

One should first note that if even R is contained in \approx , then so are \rightarrow^+ , \rightarrow^* and \leftrightarrow^* (as it is an equivalence relation), so we will focus on that case. Generally $a \rightarrow^* b$ can be seen as a way to *compute* the \approx relation, as it is directed, in a way to constrain our search space.

However, in general, we cannot deduce the converse, so it may be the case that $a \approx b$ but neither $a \rightarrow^* b$ nor $b \rightarrow^* a$ nor even is there a single c such that $a \rightarrow^* c^* \leftarrow b$, as the number of “left-right alternations” may be arbitrarily large.

The following properties are going to be very useful to deduce exactly such a converse.

Definition 6.8. We say that R is *Church-Rosser* if whenever $a \leftrightarrow^* b$, there exists some c such that

$$a \rightarrow^* c^* \leftarrow b$$

We say that R is *confluent* if whenever $b_1 \xrightarrow{*} a \xrightarrow{*} b_2$ there exists some c such that $b_1 \xrightarrow{*} c \xrightarrow{*} b_2$.

We say that R is *locally confluent* if whenever $b_1 \leftarrow a \rightarrow b_2$ there exists some c such that $b_1 \xrightarrow{*} c \xrightarrow{*} b_2$.

We say that (an arbitrary) a is in *normal form* (or a is a normal form) if there is no $a' \neq a$ such that $a \rightarrow a'$, and that R is *weakly normalizing* if for every a , there is some a' such that $a \xrightarrow{*} a'$ and a' is in normal form.

We say that R is *strongly normalizing* if there are no infinite rewrite sequences $a_1 \rightarrow a_2 \rightarrow \dots$. In particular, a strongly normalizing R is also weakly normalizing.

It turns out that if R is strongly normalizing and Church-Rosser, and effective (we can “compute” with it) then the problem of equivalence is decidable! This is because of the following lemma.

Lemma 6.9. *If R is Church-Rosser, then any normal form is unique.*

This means that, in this situation, a and b reduce to an identical normal form c if and only if $a \leftrightarrow^* b$! This means that we have the following algorithm to decide $a \leftrightarrow^* b$ (and therefore $a \approx b$ if these relations coincide):

- (1) Repeatedly apply R to a and b until normal forms a' and b' are found for them (this is possible because R is strongly normalizing).
- (2) Compare a' and b' for exact equality (sometimes called “syntactic equality”).
- (3) If $a' = b'$, we can conclude $a \leftrightarrow^* b$.
- (4) If $a' \neq b'$ we can conclude that they are *not* equivalent due to the lemma.

Note that weak normalization does not change much here except at step 1, where we need to pick reductions which eventually bring the elements to normal forms.

The strategy is therefore, for a given \approx to find an R which is (strongly) normalizing and Church-Rosser, and such that $\leftrightarrow^* = \approx$. This is roughly the goal of the entire field of *completion*. We call such an R *complete for \approx* .

The task is helped by the following facts, which we state here also without proof.

Theorem 6.10.

- (1) R is Church-Rosser iff it is confluent.
- (2) (Newman’s lemma) if R is strongly normalizing, then R is confluent iff it is locally confluent.

The example above is both locally confluent and normalizing, and therefore confluent.

This strategy can be leveraged by looking at the particulars of the equivalence relation of interest, namely quantified equations over syntactic trees, and a theory Γ , which we will usually take to be finite (usually it will have a single equation!).

We will therefore consider relations over the set of elements of the free magma M_X , and the aim is to find a rewrite system R is complete for \simeq .

Certainly \simeq is closed over substitutions, and be a *congruence*: if $a \simeq a'$ and $b \simeq b'$ under Γ , then $a \diamond b \simeq a' \diamond b'$ under Γ as well.

We therefore consider R to be both closed under substitutions and a congruence. A convenient way to represent this is via a *rewrite system*: simply a set of pairs of words $(l, r) \in M_X$ (we typically write $l \rightarrow r$) which represents the smallest congruence, closed by substitutions that contains those pairs.

Naturally, a set of laws $w \simeq w'$ can be seen, given a choice of orientation (left-to-right or right-to-left) for each law as such a rewrite system. In this case, it is very clear that the reflexive transitive closure \leftrightarrow^* recovers the original equational theory $\Gamma \models \cdot \simeq \cdot$. However, in many cases systems will either be not strongly normalizing, or confluent, or both.

For example, it is clear that commutativity (the rule $x \diamond y \simeq y \diamond x$) cannot possibly be oriented in such a way as to be terminating. Here is a non-confluent example:

$$x \diamond (y \diamond z) \rightarrow y$$

We have $a \diamond (b \diamond (c \diamond d)) \rightarrow^* b$, but also $a \diamond (b \diamond (c \diamond d)) \rightarrow^* a \cdot c$ for any a, b, c, d (which are all in normal form).

Knuth and Bendix [20] described a technique by which a theory or set of equations Γ could be turned into a complete system. The crucial idea is the observation that the non-local confluence of a rewrite system can be reduced to a finite (if the system is finite) set of “worst offenders” for confluence. If these pairs can be *joined* (reduced to the same term) then the system is confluent. It is possible to compute such pairs.

The high-level idea is therefore to identify such pairs, and add them as an unoriented equation, to be oriented if possible, and repeating until no un-joinable pairs exist. If this procedure succeeds and terminates, the system is successfully completed, and as a result the theory Γ is decidable, via the completed system as described above.

The example given above (6.7) is exactly the Knuth–Bendix procedure applied to the single equation $x \diamond (y \diamond x) = y$ (which is Equation 14 in our numbering); and indeed, it is possible to decide any word problem over (free) magmas satisfying this equation by repeatedly applying those rules, and examining the normal forms for syntactic equality.

We use the intuitive notions of “position in a word” and “word at a position p ”. We denote by $w[w']_p$ the word w with w' inserted at position p .

Definition 6.11. Given a rewrite system R and two rules $\rho_1: l_1 \rightarrow r_1$ and $\rho_2: l_2 \rightarrow r_2$ in R , we say that (t, u) is a *critical pair* for ρ_1 and ρ_2 if there is some non-variable position p in l_1

such that l_2 unifies with the term at that position. We denote by σ the most general unifier thus obtained and have $t = r_1\sigma$ and $u = l_1\sigma[r_2\sigma]_p$.

Note that, in the above setting, $t \leftarrow l_1\sigma \rightarrow u$, giving us a candidate for non-local-confluence. The next lemma states that these candidates are the most general ones.

Theorem 6.12. *Given a rewrite system R , if for every critical pair (t, u) of R , there is a term v such that $t \rightarrow^* v \leftarrow^* u$, then R is locally confluent.*

Note that building critical pairs of a finite system is computable. Therefore the only step of the completion process which require genuine creativity is the choice of the orientation of the equations, along with the proof that that orientation is strongly normalizing.

We note that even in the event that such an orientation is not found, one can still partially apply the completion procedure, using any well-founded order on terms that is stable by substitution and congruence, to obtain a semi-decision procedure for equality. This process is sometimes called *unfailing completion* and is at the core of the *superposition calculus* used in Vampire.

The associative and commutative laws, $x \diamond (y \diamond z) \simeq (x \diamond y) \diamond z$ and $x \diamond y \simeq y \diamond x$ both have confluent rewrite systems associated to them (by orienting the equations left-to-right). But it is easy to see that only the former has a complete rewrite system associated to it.

As another example, $x \diamond (x \diamond x) \simeq x$ has a natural orientation that is normalizing and such that the critical pairs are joinable, and hence is complete. This approach works for many rules.

It is easy to define a matching invariant for a complete system R associated to Γ , by simply taking as model the elements of \mathcal{M}_X , the free magma without any relations, restricted to the elements $t \in \mathcal{M}_X$ that are irreducible under R . The interpretation function ϕ simply takes a term to its normal form (under R). The lemma above ensures that this is preserved under Γ .

6.4. Unique factorization. In general, the free magma $\mathcal{M}_{X,E}$ for a given equational law E , which we can canonically define as \mathcal{M}_X / \sim_E , is hard to describe explicitly; indeed, from the undecidability of implications between equational laws, such a magma cannot be computably described for arbitrary E . Nevertheless, for some laws it is possible to obtain some partial understanding of $\mathcal{M}_{X,E}$ from a syntactic perspective. For instance, if we can refute the equivalence $w'_1 \sim_E w'_2$ by constructing a counterexample magma M that obeys E but not $w'_1 \simeq w'_2$, then this implies that the representatives $\iota_{X,E}(w'_1), \iota_{X,E}(w'_2)$ of w'_1, w'_2 in $\mathcal{M}_{X,E}$ are distinct.

We illustrate this approach with equations E of the left-absorptive form

$$(10) \quad x \simeq x \diamond f(x, y, z)$$

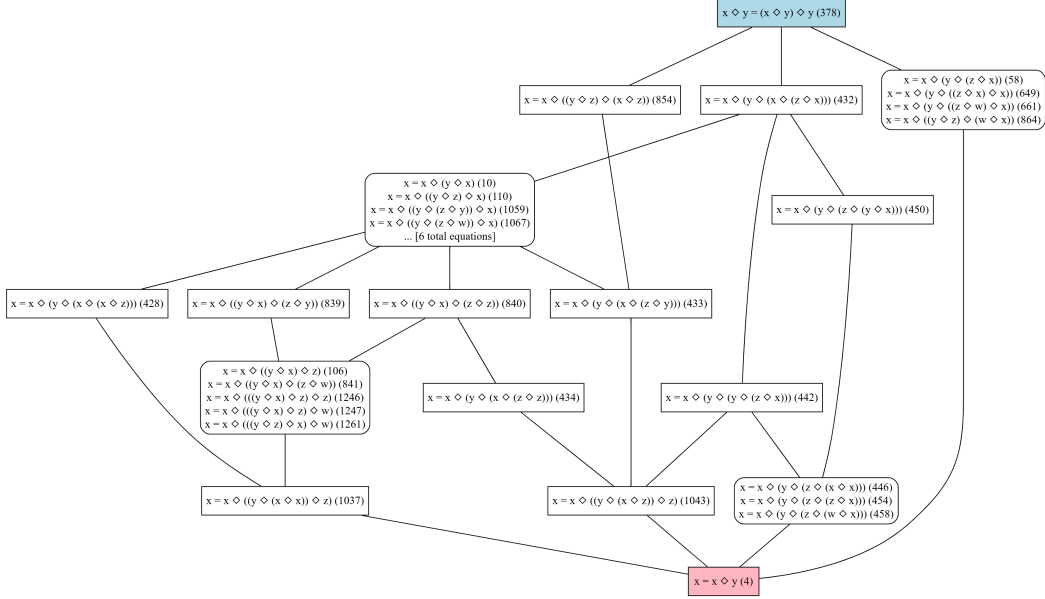


FIGURE 7. Equations similar to (E854) that are of the form Equation (10) (possibly involving a fourth indeterminate w) and imply (E378). For brevity, 70 equations equivalent to (E4) have been omitted.

for some word $f(x, y, z)$, that are also known to imply the right-idempotent law (E378). An illustrative example is the law (E854) depicted in Figure 1. Other examples are listed in Figure 7.

Lemma 6.13. *Equation (E854) is of the form (10) and implies (E378).*

Proof. Clearly we have (10) with $f(x, y, z) := (y \diamond z) \diamond (x \diamond z)$. From 10 we have in any magma obeying (E854) that

$$x = x \diamond f(x, S^2x, x) = x \diamond S(x \diamond S^2x) = x \diamond S(x \diamond f(x, x, x)) = x \diamond Sx.$$

This implies from a further application of Equation (10) that

$$y = y \diamond f(y, x, y) = (y \diamond Sy) \diamond ((x \diamond y) \diamond Sy) = f(x \diamond y, y, Sy)$$

and hence by Equation (10) again

$$(x \diamond y) \diamond y = x \diamond y$$

giving (E378). □

Let E be a law of the form Equation (10) that implies (E378). We define a directed graph \rightarrow_E on words in \mathcal{M}_X by declaring $w' \rightarrow_E w$ if $w \sim_E w'' \diamond w'$ for some $w'' \in \mathcal{M}_X$. By (E378) (applied to the quotient magma $\mathcal{M}_{X,E} = \mathcal{M}_X / \sim_E$), this is equivalent to requiring that $w \sim_E w \diamond w'$. In particular, from Equation (10) we have $f(x, y, z) \rightarrow x$ for all x, y, z . Furthermore, the relation \rightarrow_E factors through \sim_E : if $w \sim_E \tilde{w}$ and $w' \sim_E \tilde{w}'$, then $w' \rightarrow_E w$ if and only if $\tilde{w}' \rightarrow_E \tilde{w}$.

Call a word $w \in M_X$ *irreducible* if it is not of the form $w = w_1 \diamond w_2$ with $w_2 \rightarrow_E w_1$. We can partially understand the equivalence relation \sim_E on irreducible words:

Theorem 6.14 (Description of equivalence). *Let E be an equation of the form (10). Let w be an irreducible word, and let w' be a word with $w \sim_E w'$.*

(i) *If w is a product $w = w_1 \diamond w_2$, then w' takes the form*

$$w' = (((w'_1 \diamond w'_2) \diamond v_1) \diamond \dots \diamond v_n)$$

for some $w'_1 \sim_E w_1$, $w'_2 \sim_E w_2$, some $n \geq 0$, and some words v_1, \dots, v_n such that for all $0 \leq i < n$, v_{i+1} is of the form

$$v_{i+1} \sim_E f(x_i, y_i, z_i)$$

for some x_i, y_i, z_i with

$$x_i \sim_E (((w'_1 \diamond w'_2) \diamond v_1) \diamond \dots \diamond v_i).$$

In particular, $v_{i+1} \rightarrow_E x_i$.

(ii) *Similarly, if $w \in X$ is a generator of M_X , then w' takes the form*

$$w' = ((w \diamond v_1) \diamond \dots \diamond v_n)$$

for some $n \geq 0$, and some words v_1, \dots, v_n such that for all $0 \leq i < n$, v_{i+1} is of the form

$$v_{i+1} \sim_E f(x_i, y_i, z_i)$$

for some x_i, y_i, z_i with

$$x_i \sim_E ((w \diamond v_1) \diamond \dots \diamond v_i).$$

In particular, $v_{i+1} \rightarrow_E x_i$.

Conversely, any word of the above forms is equivalent to w .

Proof. We just verify claim (i), as claim (ii) is similar. The converse direction is clear from (10) (after quotienting by \sim_E), so it suffices to prove the forward claim. By the Birkhoff completeness theorem, it suffices to prove that the class of words described by (i) is preserved by any term rewriting operation, in which a term in the word is replaced by an equivalent term using (10). Clearly the term being rewritten is in w'_1 or w'_2 then the form of the word is preserved, and similarly if the term being rewritten is in one of the v_i . The only remaining case is if we are rewriting a term of the form

$$x_i = (((w'_1 \diamond w'_2) \diamond v_1) \diamond \dots \diamond v_i).$$

If $i > 0$ we can rewrite this term down to x_{i-1} , and this still preserves the required form (decrementing n by one). If $i = 0$ then we cannot perform such a rewriting because of the irreducibility of $w_1 \diamond w_2$ and hence $w'_1 \diamond w'_2$. Finally, we can rewrite x_i to $x_i \diamond v$ where v is of the form

$$v_i = f(x_i, y, z),$$

and after some relabeling we are again of the required form (now incrementing n by one). This covers all possible term rewriting operations, giving the claim. \square

Specializing to the case where w, w' are both irreducible, we conclude

Corollary 6.15 (Unique factorization). *Two irreducible words w, w' are equivalent if and only if they are either the same generator of X , or are of the form $w = w_1 \diamond w_2$, $w' = w'_1 \diamond w'_2$ with $w_1 \sim_E w'_1$ and $w_2 \sim_E w'_2$.*

As an application of this corollary, we establish

Proposition 6.16 (E854 does not imply E3316). *Equation (E854) does not imply (E3316).*

Proof. (Sketch) We work in the free group \mathcal{M}_X on two generators $X = \{x, y\}$. It suffices to show that

$$x \diamond y \not\sim_{E854} x \diamond (y \diamond (x \diamond y)).$$

Suppose this were not the case, then by Theorem 6.15 one of the following statements must hold:

- (i) $y \rightarrow_{E854} x$.
- (ii) $(y \diamond (x \diamond y)) \rightarrow_{E854} x$.
- (iii) $y \diamond (x \diamond y) \sim_{E854} y$.

If (i) holds, then we have $x \diamond y = x$ must hold in \mathcal{M}_X / \sim_E , hence (E854) would imply (E4). However, it is possible to refute this implication by a finite counterexample.

Similarly, if (iii) held, then (E854) would have to imply (E10), but this can also be refuted by a finite magma.

Finally, if (ii) held, then the claim

$$x \diamond y \sim x \diamond (y \diamond (x \diamond y))$$

to refute simplifies to

$$x \diamond y \sim x$$

and we are back to (i), which we already know not to be the case. □

7. PROOF AUTOMATION

In this project we used proof automation in two ways: automated theorem provers (ATPs) and Lean tactics. ATPs are generally stand-alone tools that implement a (semi-)decision procedure for a given formal language or related set of languages. For example, Vampire [23] is an ATP focused primarily on first-order logic using superposition, which we used extensively in this project.

ATPs are complex software that can contain bugs. Instead of trusting ATP output, we used proof certificates, which many ATPs can produce, to reconstruct proofs in Lean. The details of proof reconstruction depend on the form of the proof certificate produced by the ATP. We expand on this in Section 7.2.

Tactics in Lean, on the other hand, are meta-programs [10] that build proofs. In other words, they essentially take Lean code as input and produce Lean code as output. In this manner,

they look like another keyword in the language, and are tightly integrated by producing proofs directly. Under the hood, their implementation can be arbitrarily complex, from syntactic sugar to full decision procedures. The `duper` tactic [6], for example, implements a superposition calculus, similar to Vampire’s, but for dependent types — Lean’s underlying logical foundation.

In the rest of this section we describe the different proof automation techniques used in this project. We first discuss the different proof methods used: primarily superposition and equational reasoning, we then discuss the integration in Lean, and finally we report some basic empirical results from this project.

7.1. Proof Techniques. The main two families of ATPs and tactics we used are based on superposition/saturation and equational reasoning. In this context we also include SMT solvers, which combine specific decision procedures for theories, like congruence closure for equational reasoning, with satisfiability (SAT) solving [7]. Finally, we also used `aesop` [27], which implements a version of tableau search. This was used mainly to help specific constructions in refutations, and is not specific to proving or disproving magma implications in this sense. We describe our use of `aesop` in Section 7.2 below.

Saturation. Most of the ATPs used extensively in this project rely primarily on saturation procedures in the superposition calculus. For example, this is the case for Vampire [23].¹⁰ The core idea of these provers is that they take a set of assumptions and a conjecture, expressed in — say — first-order logic. The conjecture is negated and added to the set of assumptions, which are all put into a normal form. The ATP then tries to refute the negation by applying rules of an underlying calculus, until a proof of false (a contradiction) is derived. In this case, the conjecture was (classically) true, and the ATP has found a proof by contradiction, often called a “refutation” or “saturation” proof.

The underlying calculi vary from system to system, but they often have a variant of a resolution clause of the form:

$$\frac{C \vee L \quad D \vee \neg L}{C \vee D}$$

This can be read as $C \vee L$ with $D \vee \neg L$ implies $C \vee D$, where C, D, L are formulas in e.g. first-order logic. Superposition calculi have a variant of this rule that deals with equality directly, and thus are more efficient at reasoning about equality.

In this project we used Vampire [23], Duper [6] and Prover9 and Mace4 [28] which are all based on variants of saturation for proving.

here we could add a screenshot of using vampire or prover9

Equational Reasoning. As already discussed in Section 6.3, equational reasoning is a type of reasoning based on equational logic and rewriting with congruence [3]. In general, an equational reasoning procedure takes a series of equations and determines whether another equation can be deduced from it. A core tool in equational reasoning are e-graphs, a data

¹⁰See also [4] for a gentler exposition.

structure used to represent congruence classes of terms. By themselves, e-graphs provide an efficient means of implementing a decision procedure for congruence closure over ground equations (i.e. equations without variables). Extensions to this procedure, for example by quantifier instantiation via e-matching [8], also allow for a semi-decision procedure for congruence closure over non-ground equations.

SMT solvers like Z3 [9] use equational reasoning for deciding the theory of equality with uninterpreted functions [24, 8]. On the other hand, equality saturation [46] uses e-graphs by extending congruence closure to a more controlled search, enabling optimization and conditional rewriting. One of the main advantages of using equational reasoning to reason about implications of magma laws is that we get very explicit proofs: a proof that $l \vdash l'$ is given by a sequence of rewrites that starts at the left-hand side of l' and arrives at the right-hand side through applications of l .

In this project we used Z3 [9], Prover9 and Mace4 [28], a custom ATP for magmas based on egg [46], and the Lean `egg` tactic [21, 38, 39], which all work with equational logic. We have also reasoned with manual (custom written) heuristics about simple rewrites.

7.2. Proof Reconstruction and Integration.

Improve the structure of this section

While ATPs are very useful for solving theorems in this project, they don't integrate with Lean out of the box. ATPs may potentially produce unsound proofs, or worse, derive incorrect results. Thus, by default, theorems in Lean cannot be proven by deferring to the result of an ATP. Instead, the results of an ATP can be used to (re-)construct a proof of the form required by Lean. There are different approaches to proof (re-)construction which we employ in this project.

An exception for this are tactics like `duper`, `aesop` and `egg`. Figure 8 shows an example of the `egg` tactic as used in this project. It integrates directly in Lean, generating a Lean proof directly. With the variant `egg?`, depicted in the screenshot, it uses an auxiliary tactic `calcify`¹¹ to generate a human-readable proof as a series of calculation steps, which can be incorporated into the file with a single click.

In general, integration implies two steps: invoking the decision procedures/ATPs (translating the problem from Lean into the languages and logics they use), and conversely, (re-)constructing the results from the decision procedure as a (persistent) Lean proof. These two aspects present different challenges, and require different strategies, depending mostly on the kind of proof strategy the decision procedure uses.

For saturation proofs ...

add explanations from the following URL

¹¹<https://github.com/nomeata/lean-calcify>

```

graph.lean •
  Subgraph.lean > { } Subgraph > Equation14_implies_Equation23
  namespace Subgraph
  set_option egg.explosion true
  /- Obtained with lean-egg -/
  @[equational_result]
  theorem Equation14_implies_Equation23 (G: Type*) [Magma G] (h: Equation14 G) : Equation23 G :=
  by egg? [*] Try this: ⚡ intro x ⚡ calc ⚡ x ⚡ _ = (x ⚡ x) ⚡ (x ⚡ (x ⚡ x)) := (h x (x ⚡ x)) ⚡ _ =
  Try this:
  @[equ]
  theor
  calc
  | by
  x
  _ = (x ⚡ x) ⚡ (x ⚡ (x ⚡ x)) := (h x (x ⚡ x))
  @[equ]
  theor
  _ = (x ⚡ x) ⚡ x := congrArg (fun a' ↦ (x ⚡ x) ⚡ a') (Eq.symm (h
  x x)) Lean 4
  Equation14 G :=
  by
  A sequence of tactics in brackets, or a delimiter-free indented sequence of tactics.
  Delimiter-free indentation is determined by the first tactic of the sequence.
  @[equ]
  theor
  View Problem (Alt+F8) Quick Fix... (Ctrl+)
  by egg [*]
  Equation381 G :=
  
```

FIGURE 8. An example of the use of proof automation with tactics. This shows the egg tactic as it was used to generate (human-readable) equational proofs of positive implications.

<https://leanprover.zulipchat.com/#narrow/channel/458659-Equational/topic/Vampire.20-.3E.20Lean/near/478147138>

For equational proofs from external provers (e.g. MagmaEgg), we also used a simple version of reconstruction by (re-)constructing the proofs of equality from an explanation, using congruence lemmas from Lean. In the case of equational proofs from the egg tactic, these could be converted into a series of calc steps, documenting the explicit calculation, using the calcify tactic.

In general, we have observed that there are multiple ways of integrating decisions procedures within Lean, with different levels of integration.

- (1) Using a Lean tactic, which calls a decision procedure written in Lean (like aesop or duper).
- (2) Using a Lean tactic, which calls an existing (external) ATP and reconstructs a proof term from the ATP’s result (like bv_decide or egg).
- (3) Using an external script which calls an existing ATP and generates a source file .lean which captures the result explicitly.

This project primarily used the least integrated approach, (Option 3), as it was fastest and required no dependencies on the other contributors. This also has drawbacks: primarily, while the upfront effort is lower, the effort to use is higher than with a tactic, once the tactic is developed. It also makes integration with larger code bases more difficult: in this project the (mathematical) dependencies were by design very minimal, magmas are simple, and we built our own definitions for them. For example, integration with the typeclass system becomes much more difficult when working with more complex mathematical objects that build on multiple, nested layers of structure in non-trivial ways. In that case, for example, tactics need to synthesize typeclass instances, deal with diamonds and different notions of equivalence [45].

Semi-Automated Counterexample Guidance. Another use of ATPs has been in a semi-automatic fashion, to find counter-examples. The general strategy was to use ATPs to find counter-examples to implications by building magmas iteratively. If we want to build a counterexample to $l \vdash l'$, we want to construct a magma where l holds but l' does not. In this method, we iteratively strengthen a construction with additional hypotheses, and use the ATP to check whether these hypotheses are not too strong (to imply l') or unsound (to disallow l).

this should also be expanded more, at least with references to some of the constructions in other chapters.

While equational reasoning can also be used in a semi-automatic fashion to prove equations [21], the positive implications in the main implication graph of project were all simple enough that we did not need a semi-automatic approach for them.

discuss guided search in the finite implications or the Higman-Neumann work Jose Brox has done.

maybe add a screenshot here of the workflow of using a seed to find counterexamples with prover9 or vampire?

7.3. Empirical Results. Finally, we report some empirical results from use of ATPs for this project, in terms of performance. The aim of this section is not to be a careful evaluation and benchmark comparison of the different ATPs; instead, we present our work here as a more informal “field report” documenting our experiences. In particular, we do not draw firm conclusions about the overall capabilities¹² of the different ATPs. Rather, this serves as a use-case documenting the experience of (mostly) novice users.

throw a couple of "benchmarking" tables for the same ATP with different parameters and for different ATPs, talk about some relative gains in time (changing parameters we saw a 500 times speedup on this particular problem), etc. This is knowledge I think we have gained to some extent, and certainly I would have been glad to receive this kind of hints before we started!". Then leave it as an interesting open problem to properly develop and measure benchmarks for ATPs based on this project.

Any comparative study of semi-automated methods with fully automated ones? In principle, the semi-automated approach could be more automated using a script or "agent" to call various theorem provers. See this discussion

¹²One reason for this is that different ATPs were deployed at different stages of the project. In particular, the later ATP runs were performed in an environment when a large fraction of the implications had already been settled, and only a small remainder set was tested by our project. The current set of ATP-generated formalized implications has also been subject to a number of reductions to optimize compilation time, so we would caution against reading too much into the raw number of such formalizations in the codebase.

See this discussion on the value of using different ATPs and setting run time parameters etc. at different values.

What are the hardest implications to prove? See this discussion.

8. IMPLICATIONS FOR FINITE MAGMAS

Expand this sketch

Recap discussion from <https://leanprover.zulipchat.com/#narrow/channel/458659-Equational/topic/Austin.20pairs>

9. HIGMAN–NEUMANN LAWS

report on Higman–Neumann laws

10. AI AND MACHINE LEARNING CONTRIBUTIONS

As discussed in Section 7, the ETP made extensive use of automated theorem provers in completing the primary goal of determining and then formalizing all the implications between the specified equational laws. In contrast, we were only able to utilize modern large language models (LLMs) in a fairly limited fashion. Such models were useful in writing initial code for the graphical user interfaces discussed in Section 11, as well as performing some code autocompletion (using tools such as *Github Copilot*) when formalizing an informal proof in Lean. In one instance, *ChatGPT* was used¹³ to guess a complete rewriting system for the law $x \diamond ((y \diamond y) \diamond z) \simeq x \diamond y$ (E1659) which could then be formally verified, thus resolving all implications from this equation. However, in most of the difficult implications that resisted automated approaches, we found that LLMs did not provide useful suggestions beyond what the human participants could already propose.

On the other hand, we found that machine learning (ML) methods showed some promise of being able to heuristically predict the truth value of portions of implication graph, as we shall now discuss¹⁴.

10.1. Graph ML: Directed link prediction on the implication graph. We experimented with various Graph Neural Network (GNN) autoencoders to predict missing edges, providing a way to estimate the truth value of unproven implications. To assess these models, we defined three test sets focusing on edge existence, directionality, and bidirectionality.

¹³<https://chatgpt.com/share/670ce7db-8a44-800d-a5dc-8462c12eca3b>

¹⁴For some discussion of other ML experiments performed during the project, see <https://leanprover.zulipchat.com/#narrow/channel/458659-Equational/topic/Machine.20learning.2C.20first.20results>.

The results give insight into how these models handle dense, directed graphs like ours. A more detailed report follows below.

10.1.1. *Motivation.* Directed Link Prediction [15] is a method enabling machine learning and deep learning models to predict missing edges in a directed graph. For our implication graph, this translates to predicting the truth values of unproven implications. This task serves as a necessary first step for advancing in the following directions:

- (1) **Reasoning over mathematical knowledge graphs:** Recent advancements allow language models to integrate information from multiple modalities. For example, [49, 48] share information between corresponding layers of Language Models (LMs) and Graph Neural Networks (GNNs), enabling simultaneous learning from text corpora and graph-structured data within the same expert domain. By leveraging both modalities, the language model can better *structure* its knowledge and respond to complex queries. This dual learning process combines masked language modeling for text with link prediction on the graph, highlighting the importance of link prediction for robust reasoning.
- (2) **Higher-order implication graphs:** Our implication graph currently represents only implications of the form $p \implies q$, not more complex ones like $(p \wedge r) \implies q$. Extending to such higher-order edges would likely involve connecting sets of nodes, thereby requiring hypergraph representations. For a systematic overview, see [18]. While specific hypergraph neural architectures exist [12], we believe it is still conceptually important to apply Directed Link Prediction to simpler implication graphs first, providing insights and guiding principles that can anticipate challenges in higher-order graph representation learning.

10.1.2. *Data.* We used this edge list, generated on October 20, 2024, with the following commands:

```
lake exe extract_implications outcomes > data/tmp/outcomes.json
scripts/generate_edgelist_csv.py
```

The structure of the implication graph is summarized below:

```
graph_summary = {
  "total_nodes": 4694,
  "total_directed_edges": 8178279,
  "edge_density_percentage": 37, # Percentage of possible edges that exist
  "bidirectional_edges": 2475610,
  "bidirectional_percentage": 30 # Percentage of all edges that are bidirectional
}
```

Below is a summary of the edge types in the graph:

```
edge_counts = {
```

```

    "explicit_conjecture_false": 92,
    "explicit_proof_false": 582316,
    "explicit_proof_true": 10657,
    "implicit_conjecture_false": 142,
    "implicit_proof_false": 13272681,
    "implicit_proof_true": 8167622,
    "unknown": 126
  }

```

Edges are labeled according to the following scheme:

```

edge_labels = {
  "implicit_proof_true": 1,
  "explicit_proof_false": 0,
  "implicit_proof_false": 0,
  "explicit_proof_true": 1,
  "explicit_conjecture_false": 0,
  "implicit_conjecture_false": 0,
  "unknown": 0
}

```

The `unknown` class contains a very small number of edges (126), so their impact on the training phase is expected to be negligible. Future approaches might address this class by excluding `unknown` edges from the training set.

10.1.3. *Methods.* Consider a directed graph $G = (V, E)$ where $E = \{(u, v) \mid u, v \in V\}$ is the edge set and $|V| = n$. We assume that each node is associated with a feature vector, resulting in an $X \in \mathbb{R}^{n \times f}$ feature matrix.

We define the existing edges $(a, b) \in E$ as *positives* and the non-existing edges $(c, d) \notin E$ as *negatives*.

Intuitively, performing Directed Link Prediction (DLP) on G involves randomly splitting E into three disjoint sets: E_{train} , E_{val} , and E_{test} , such that:

- E_{train} is the training set,
- E_{val} is the validation set,
- E_{test} is the test set,
- and $E = E_{\text{train}} \dot{\cup} E_{\text{val}} \dot{\cup} E_{\text{test}}$.

The model then learns from $G_{\text{train}} = (V, E_{\text{train}})$ to map the topological and feature-related information of two nodes u and v to a probability p_{uv} that $(u, v) \in E_{\text{test}}$.

However, this setup presents two key issues among others:

- (1) The model learns only from *positives*, so it cannot recognize *negatives*.

- (2) The model is evaluated only on *positives*, preventing us from measuring its ability to identify *negatives*.

To address these limitations, we adopted the setup proposed in [40]. Specifically, we redefined E_{train} to E_{train}^p (*positives*) and introduced:

$$E_{\text{train}} = E_{\text{train}}^p \dot{\cup} E_{\text{train}}^n$$

where E_{train}^n includes all possible *negatives* in $G_{\text{train}} = (V, E_{\text{train}}^p)$. The model is now required to predict the non-existence of edges in E_{train}^n .

Similarly, if we redefine E_{test} as follows:

$$E_{\text{test}} = E_{\text{test}}^p \dot{\cup} E_{\text{test}}^n$$

where E_{test}^n is a *random* sample of *negatives*, the model’s evaluation would fail to capture two crucial aspects:

- (1) The model’s ability to distinguish (u, v) from (v, u) for all $(u, v) \in E_{\text{test}}^p$.
- (2) The model’s ability to identify bi-implications.

These limitations arise from the random selection of negative edges in E_{test}^n . To address this, we define three distinct test sets: E_{test}^G , E_{test}^D , and E_{test}^B , to evaluate different facets of the model’s performance:

- **General Test Set** (E_{test}^G): Here, $E_{\text{test}} = E_{\text{test}}^p \dot{\cup} E_{\text{test}}^n$, where E_{test}^n is a random sample of non-existent edges with the same cardinality as E_{test}^p . This set assesses the model’s ability to detect the presence of edges, regardless of direction. A model that cannot distinguish edge direction may still perform well on this set, highlighting the need for the following two additional test sets.
- **Directional Test Set** (E_{test}^D): Defined as $E_{\text{test}}^{\text{up}} \dot{\cup} \tilde{E}_{\text{test}}^{\text{up}}$, where $E_{\text{test}}^{\text{up}}$ consists of unidirectional edges in E_{test}^p , and $\tilde{E}_{\text{test}}^{\text{up}}$ contains their reverses (negatives by construction). This set evaluates the model’s ability to recognize edge direction, making it suitable for assessing direction-sensitive models.
- **Bidirectional Test Set** (E_{test}^B): Defined as $E_{\text{test}}^{\text{bp}} \dot{\cup} E_{\text{test}}^{\text{n}}$, where $E_{\text{test}}^{\text{bp}}$ contains one direction of all bidirectional edges in E_{test}^p , and $E_{\text{test}}^{\text{n}} \subset \tilde{E}$ includes a subset of their reverses. This set evaluates the model’s ability to identify bi-implications, as each edge in E_{test}^B has a reverse that is positive, but only half are bidirectional in practice.

We tested the following models:

- **GAE** [15]
- **Gravity-GAE** [40]
- **Source/Target-GAE** [40]

- **DiGAE** [22]
- **MagNet** [50]

All these models are graph-based autoencoders with distinct encoder-decoder architectures. Notably, GAE is the only model structurally unable to differentiate edge directions. Each model outputs the probability that an ordered pair of nodes has a directed edge between them, with nodes represented using one-hot encodings as features.

We trained the models using Binary Cross Entropy as the loss function, balancing the contribution of positive and negative edges through re-weighting. On the *General* test set, we evaluated the following metrics:

- **AUC** (Area Under the ROC Curve): Measures the probability that the model ranks a random positive edge higher than a random negative edge. Higher values indicate better discrimination between positive and negative edges.
- **AUPRC** (Area Under Precision-Recall Curve): Assesses model performance, particularly in cases of class imbalance. AUPRC is more robust to imbalanced data than AUC.
- **Hits@K**: Evaluates the fraction of times a positive edge is ranked within the top K positions among personalized negative samples [26]. Briefly, given a positive edge, its M personalized negative samples are M negative edges with the same head but different tails. We calculate Hits@K for $K = 1, 3, 10$ to assess ranking quality, and set $M = 100$. Therefore, Hits@K = 0.1 means that on average, the model ranks a positive edge within the highest-ranked K personalized negatives 10% of the time.
- **MRR** (Mean Reciprocal Rank): Computes the average reciprocal rank of positive edges among their personalized negative samples [26] (the same as those used for Hits@K) providing an overall measure of ranking accuracy. For instance, $MRR = 0.1$ means that on average, the model ranks a positive edge as 10th among M personalized negative samples.

Each metric ranges from 0 to 1, with higher values reflecting improved performance. Based on prior work, we expect AUC and AUPRC scores to approach 1, while MRR and Hits@K often yield values around 0.15 for similar undirected tasks [26]. *Directional* and *Bidirectional* performances were evaluated using only AUC and AUPRC, since Hits@K and MRR are hardly definable in those scenarios. The number of training epochs was optimized through Early Stopping on the *General* validation set performance (given by the sum of AUC and AUPRC).

10.1.4. *Results.* The results below represent average performance over six random splits of E_{train} , E_{val} , and E_{test} while keeping the model’s seed fixed for fair comparison. The *training* / *validation* / *test* split proportions are:

- 85/5/10 for unidirectional edges,
- 65/15/30 for bidirectional edges.

| Model | G ROC AUC | G AUPRC | G Hits@1 | G Hits@3 |
|------------------|--------------------------------|---------------------------------|---|---|
| gae | $0.8484 \pm 9 \times 10^{-4}$ | $0.8558 \pm 6 \times 10^{-4}$ | $6 \times 10^{-5} \pm 4 \times 10^{-5}$ | $6 \times 10^{-5} \pm 4 \times 10^{-5}$ |
| gravity_gae | $0.9806 \pm 3 \times 10^{-4}$ | $0.9753 \pm 4 \times 10^{-4}$ | $0.069 \pm 6 \times 10^{-3}$ | $0.101 \pm 5 \times 10^{-3}$ |
| sourcetarget_gae | $0.99976 \pm 1 \times 10^{-5}$ | $0.999736 \pm 8 \times 10^{-6}$ | $0.077 \pm 4 \times 10^{-3}$ | $0.147 \pm 7 \times 10^{-3}$ |
| mlp_gae | $0.99315 \pm 1 \times 10^{-5}$ | $0.99409 \pm 1 \times 10^{-5}$ | $0.181 \pm 7 \times 10^{-3}$ | $0.299 \pm 7 \times 10^{-3}$ |
| digae | $0.9978 \pm 3 \times 10^{-4}$ | $0.998 \pm 3 \times 10^{-4}$ | $0.035 \pm 6 \times 10^{-3}$ | $0.068 \pm 1 \times 10^{-2}$ |
| magnet | $0.989 \pm 1 \times 10^{-4}$ | $0.99076 \pm 3 \times 10^{-5}$ | $0.151 \pm 1 \times 10^{-2}$ | $0.26 \pm 2 \times 10^{-2}$ |

TABLE 2. Results for various graph autoencoder models.

10.1.5. *Discussion.* We observe consistently high *General* AUC and AUPRC scores, close to 1. These high values are expected, as similar neural architectures have demonstrated strong performance in graphs of comparable size [15]. The high ratio of existing to non-existing edges in the implication graph (approximately 37%) likely contributes to the near-perfect *General* AUC and AUPRC scores. For context, benchmark datasets such as Cora and Citeseer (e.g., directed and undirected) contain fewer than 1% of all possible edges.

Interestingly, the GAE model, though structurally unable to distinguish edge direction, performs well on the *General* task (if we consider AUC and AUPRC only). This experimentally confirms the need for including *Directional* and *Bidirectional* test sets, which allow comprehensive evaluation across all facets of Directed Link Prediction (DLP).

All other models demonstrate high AUC and AUPRC scores across the *General*, *Directional*, and *Bidirectional* test sets, indicating strong predictive capabilities even when accounting for directionality and bidirectionality.

Notably, the `mlp_gae` and `magnet` models also achieve competitive scores in MRR and Hits@K metrics.

10.1.6. *Conclusions.* We evaluated the performance of six different graph autoencoders on a Directed Link Prediction (DLP) task. By adopting a multi-task evaluation framework, we assessed the models comprehensively across various aspects of DLP. All models demonstrated strong performance on AUC and AUPRC metrics, with some also achieving high scores on MRR and Hits@K.

Node features were represented using one-hot encodings, meaning that the models received no explicit information about the equations represented by the nodes. Instead, they relied entirely on the topological structure encoded during training. This approach aligns with previous research suggesting that one-hot encodings can promote asymmetric embeddings [40]. However, future experiments could explore alternative representations, such as encoding the equations with text-based embeddings like Word2Vec, to potentially enhance the models' understanding of the nodes' semantic content.

In summary, our findings highlight the adaptability and robustness of graph autoencoders for DLP tasks in dense, directed graphs. This approach not only demonstrates robustness in predicting directed links but also suggests promising potential for future improvements

through enhanced feature representations, thereby advancing the capabilities of link prediction in complex mathematical knowledge graphs.

11. USER INTERFACES

A number of custom web applications were developed as part of the ETP. While many past Lean formalization projects have primarily relied on the Lean blueprint tool to organize tasks and track progress, the large volume of (transitive) implications tracked by the ETP, along with the research-oriented nature of the project, necessitated the development of custom tools to complement the blueprint tool. These web applications also made information more accessible to project participants and other interested parties, including those unfamiliar with Lean or the custom software developed for the project. The project features four primary interfaces:

- (1) The **ETP dashboard**¹⁵ displays the high-level overview of the project: the total number of resolved, conjectured, and unknown implications for the general and finite implication graphs. The dashboard also includes links to other tools, data, and visualizations about the implication graphs.
- (2) The **Equation Explorer**¹⁶ is the primary tool to navigate the implication graph. For a given equation, it displays its inbound and outbound implications, as well as other members of its equivalence class. The explorer allows navigating either the general or finite implication graphs. The explorer also features custom commentary for a given equation (when available), serving as a repository for information and links. It also links to Graphiti visualizations and an example of its smallest satisfying magma, if one exists. Figure 9 shows an example view of the explorer.
- (3) **Graphiti**¹⁷ visualizes the implication graph as a Hasse diagram, where downward edges represent subset relationships, and upward edges represent implications. Equivalence classes are collapsed into single nodes for clarity. Graphiti supports search parameters to visualize specific subsets of the graph. It can also display the entire implication graph, though the complete graph is large and challenging to navigate. Figure 7 is an example of a Graphiti visualization.
- (4) The **Finite Magma Explorer**¹⁸ tests which equations a given finite magma satisfies or fails to satisfy. Users input finite magmas as Cayley tables. The tool is aware of the finite implication graph, so if an input magma witnesses an unknown refutation, it notifies the user and provides instructions for contributing the result to the GitHub repository.

The data for these tools is extracted directly from the Lean-formalized proofs in the project's GitHub repository, ensuring it always faithfully reflects the current state of progress. Additionally, the data is automatically updated with each code change using continuous integration (CI), eliminating the need for manual updates.

¹⁵https://teorth.github.io/equational_theories/dashboard/

¹⁶https://teorth.github.io/equational_theories/implications/

¹⁷https://teorth.github.io/equational_theories/graphiti/

¹⁸https://teorth.github.io/equational_theories/fme/

Equation Details

[Back to List](#)Equation42[$x \diamond y = x \diamond z$](Dual equation: [Equation45\[\$x \diamond y = z \diamond y\$ \]](#))(Visualize [implies](#) and [implied by](#) of the equation, or see [1](#), [2](#), [3](#) graph edges away)(Size of smallest non-trivial magma: [2 \(Explore\)](#)) Hide equivalent equations Treat conjectures as unknown Display the finite graph Show only explicit proofs

This equation implies (\Rightarrow):

| Implies | Does not imply | Unknown |
|--|---|---------|
| Equation1[$x = x$] Try This! Show Proof | Equation2[$x = y$] (+ 1495 equiv.) Try This! Show Proof | None |
| Equation307[$x \diamond x = x \diamond (x \diamond x)$] Try This! Show Proof | Equation3[$x = x \diamond x$] Try This! Show Proof | |
| Equation308[$x \diamond x = x \diamond (x \diamond y)$] Try This! Show Proof | Equation4[$x = x \diamond y$] (+ 70 equiv.) Try This! Show Proof | |
| Equation309[$x \diamond x = x \diamond (y \diamond x)$] Try This! Show Proof | Equation5[$x = y \diamond x$] (+ 70 equiv.) Try This! Show Proof | |
| Equation310[$x \diamond x = x \diamond (y \diamond y)$] (+ 1 equiv.) Try This! Show Proof | Equation8[$x = x \diamond (x \diamond x)$] Try This! Show Proof | |
| Equation311[$x \diamond x = x \diamond (y \diamond z)$] Try This! Show Proof | Equation9[$x = x \diamond (x \diamond y)$] (+ 8 equiv.) Try This! Show Proof | |

This equation is implied by (\Leftarrow):

| Implied by | Not implied by | Unknown by |
|---|--|------------|
| Equation2[$x = y$] (+ 1495 equiv.) Try This! Show Proof | Equation1[$x = x$] Try This! Show Proof | None |
| Equation4[$x = x \diamond y$] (+ 70 equiv.) Try This! Show Proof | Equation3[$x = x \diamond x$] Try This! Show Proof | |
| Equation24[$x = (x \diamond x) \diamond y$] (+ 111 equiv.) Try This! Show Proof | Equation5[$x = y \diamond x$] (+ 70 equiv.) Try This! Show Proof | |
| Equation41[$x \diamond x = y \diamond z$] (+ 418 equiv.) Try This! Show Proof | Equation8[$x = x \diamond (x \diamond x)$] Try This! Show Proof | |
| Equation256[$x = ((x \diamond x) \diamond x) \diamond y$] (+ 75 equiv.) Try This! Show Proof | Equation9[$x = x \diamond (x \diamond y)$] (+ 8 equiv.) Try This! Show Proof | |
| Equation374[$x \diamond y = (x \diamond x) \diamond x$] (+ 36 equiv.) Try This! Show Proof | Equation10[$x = x \diamond (y \diamond x)$] (+ 5 equiv.) Try This! Show Proof | |

Equivalent Equations

Equivalent Equations:

- [Equation38\[\$x \diamond x = x \diamond y\$ \]](#)
- [Equation322\[\$x \diamond y = x \diamond \(x \diamond x\)\$ \]](#)
- [Equation324\[\$x \diamond y = x \diamond \(x \diamond z\)\$ \]](#)

FIGURE 9. An example of the information displayed by the Equation Explorer for a specific equation.

12. DATA MANAGEMENT

Describe how data was handled during the project and how it will be managed going forward.

All the data and formalizations generated by the project are placed in the Github repository¹⁹, while the discussion is almost entirely contained in a dedicated channel on the Lean Zulip²⁰. The implication graph can be downloaded from Equation Explorer²¹, and can also indicate the individual *Lean* theorems required to establish or refute any given implication, although currently we have only formalized a generating set of implications and refutations in *Lean*, rather than the entirety of the implication graph.

¹⁹https://github.com/teorth/equational_theories

²⁰<https://leanprover.zulipchat.com/#narrow/channel/458659-Equational>

²¹https://teorth.github.io/equational_theories/implications/

13. CONCLUSIONS AND FUTURE DIRECTIONS

This project successfully demonstrated that large-scale explorations of a space of mathematical statements - in this case, the implications or non-implications between selected equational laws - can be crowdsourced using modern collaboration platforms and proof assistants. No single tool or method was able to study the entirety of this space, and many informal proofs generated contained non-trivial errors; but there were multiple techniques that could treat significant portions of the space, and through a collaborative effort combined with the proof validation provided by *Lean*, one could synthesize these partial and fallible contributions into a complete and validated description of the entire implication graph. While this particular graph was a comparatively simple structure to analyze, we believe that this paradigm could also serve as a model for future projects devoted to exploring more sophisticated large-scale mathematical structures.

Several factors appeared to be helpful in ensuring the success of the project, including the following:

- **A clearly stated primary goal, with an end condition and precise numerical metrics to measure partial completion.** From the outset, there was a specific goal to attain, namely to completely determine and then formalize the implication graph on the original set of 4694 laws. Progress towards that goal could be measured by a number of metrics, such as the number of implications that were conjectured but unformalized, or not conjectured at all. Such metrics allowed participants to see how partial contributions, such as formalizing a certain subset of implications, advanced the project directly towards its primary goal. This is not to say that all activity was devoted solely towards this primary goal, but it did provide a coherent focus to help guide and motivate other secondary activities.
- **A highly modular project.** It was possible for any given coauthor to work on a small subset of implications and focus on a single proof technique, without needing to understand or rely upon other contributions to the project. This allowed the work to be both parallelized and decentralized; many contributors launched their own investigations broadly within the framework of the project, without needing centralized approval or coordination.
- **Low levels of required mathematical and formal prerequisites.** The problems considered in the project did not require advanced mathematical knowledge (beyond a general familiarity with abstract algebra), nor a sophisticated understanding of formal proof assistants. This permitted contributions from a broad spectrum of participants, including those without a graduate mathematical training, as well as mathematicians with no experience in proof formalization. At a technical level, it also meant that formalization of proofs into *Lean* could be done immediately once certain base definitions (such as *Magma*) were constructed. This can be compared for instance with the recent formalization of the Polynomial Freiman–Ruzsa conjecture²², in which significant effort was expended in the first few days to settle on a suitable framework to formalize the mathematics of Shannon entropy. While some more sophisticated formal structures (such as the syntactic description of laws as pairs of words in a

²²<https://github.com/teorth/pfr>

FreeMagma) were later introduced in the project, it was relatively straightforward to refactor previously written code to be compatible with these structures as they were incorporated into the project.

- **Variable levels of difficulty, and the amenability to partial progress.** Traditional mathematics projects generally involve a small number of extremely hard problems, with incomplete progress on these problems being difficult to convert into clean partial results. In contrast, the ETP studied a large number of problems with a very broad range of difficulty, so that even if a given proof strategy did not work for a given implication, it could be the case that there was some class of easier implications for which the strategy was successful. This allowed for a means to validate such ideas, and allowed the project to build up a useful and diverse toolbox of proof techniques which became increasingly necessary to handle the final and most difficult implications in the project. It also created a dynamic in which the project initially focused on easy techniques to resolve a significant fraction of the implications, gradually transitioning into more sophisticated methods that focused on a much smaller number of outstanding implications that had proven resistant (or even “immune”) to all easier approaches.
- **Centralized and standardized platforms for discussion, project management, and validation.** While the project was decentralized at the level of the participant, there was a centralized location (a channel²³ on the Lean Zulip) to discuss all aspects of the project, as well as a centralized repository²⁴ to track all contributions and outstanding issues, a centralized blueprint²⁵ to describe technical details of proofs to be formalized, and a single formal language (*Lean*) to validate all contributions. A significant portion of the activity in the early stages of the project was devoted to setting out the standards and workflows for handling both the discussion and the contributions, in particular setting up a contributions page²⁶ and adopting a code of conduct²⁷. This gave some structure and predictability to what might otherwise be a chaotic effort.
- **Development of custom visualization tools.** As discussed in Section 11, several tools were developed (in part with AI assistance) to help visualize and navigate the implication graph while it was in a partial stage of development, allowing for participants to independently identify problems to work on, and to validate and use the contributions of other participants even before they were fully formalized. For instance, a participant could propose a finite counterexample to an implication by posting a link to the magma in *Finite Magma Explorer*, allowing for immediate validation of the counterexample, or use *Equation Explorer* or *Graphiti* to observe some interesting phenomenon in the implication graph that other participants could reproduce and study.
- **Applicability of existing software tools.** As described in Section 7, many of the implications in the ETP were amenable to application of “off-the-shelf” automated theorem provers (ATPs); while some trial and error was needed to determine good

²³<https://leanprover.zulipchat.com/#narrow/channel/458659-Equational>

²⁴https://github.com/teorth/equational_theories

²⁵https://teorth.github.io/equational_theories/blueprint/

²⁶https://github.com/teorth/equational_theories/blob/main/CONTRIBUTING.md

²⁷https://github.com/teorth/equational_theories/blob/main/CODE_OF_CONDUCT.md

choices of parameters, these tools could largely be applied directly to the project without extensive customization. (However, the later transcription of ATP output into Lean was sometimes non-trivial.)

- **Receptiveness to new techniques and tools.** Crucially, the methods used to make progress on the project were not specified in advance, and contributions from participants with new ideas, techniques, or software tools that were not initially anticipated were welcomed. For instance, the theory of confluence (Section 6.3) was not initially known to the first project participants, but was brought to the attention of the project by a later contributor. Conversely, while there were hopes expressed early in the project that modern large language models (LLMs) could automatically generate many of the proofs required, it turned out in practice that other forms of automation, particularly ATPs, were significantly more effective at this task (at least if one restricted to publicly available LLMs), and the project largely moved away from the use of such LLMs (other than to help create the code for the visualization tools).

There are several mathematical and computational questions that could potentially be addressed in future work building upon the outcomes of ETP. Here is a list of some possible such future directions.

- (1) Does the law $x \simeq y \diamond (x \diamond ((y \diamond x) \diamond y))$ (E677) imply $x \simeq ((x \diamond x) \diamond x) \diamond x$ (E255) for finite magmas? This is the last remaining implication (up to duality) for finite magmas to be resolved. A number of partial results on this problem may be found at https://teorth.github.io/equational_theories/blueprint/677-chapter.html.
- (2) Does the law $x \simeq y \diamond ((y \diamond (y \diamond x)) \diamond (z \diamond y))$ (E5093) have any infinite models? In [16] it was shown that it has no finite models, but the infinite model case was left as an open question.
- (3) The ETP focused on determining relations $E \vdash E'$ between one law and another. Could the same methods also systematically determine more complex logical relations, such as $E_1 \wedge E_2 \vdash E_3$, for all laws E_1, E_2, E_3 in a specified set?
- (4) A key feature of finite magmas M is that they are surjunctive: maps from M to itself that are surjective are necessarily injective (or vice versa). Are there equational theories that admit surjunctive models (in which every definable surjective map is injective), but yet do not have any finite models?
- (5) Are there interesting examples of implications $E_1 \vdash E_2$ which are “irreducible” in the sense that there is no equational law E with $E_1 \vdash E \vdash E_2$, other than those laws equivalent to either E_1 or E_2 ?
- (6) How “stable” is a given law E ? For instance, if a finite magma obeys a law E some proportion $1 - \varepsilon$ of the time, with ε small, can the magma be perturbed into one that obeys E exactly? Related to this is the question of whether a law E is “rigid” or “mutable”: is it possible to make a small number of modifications to a magma obeying E in a way that still preserves E ? Such properties helped suggest whether certain magma construction techniques, such as modifying a base magma, were likely to be successful.

13.1. **Miscellaneous remarks.** It is possible that the timing in which certain proof methods were introduced into the project created some opportunity costs. For instance, by deploying automated theorem provers at an early stage, we might have settled some implications that had more interesting human-readable proofs that we missed. Similarly, we developed some sophisticated theory for the equation E854, such as Theorem 6.15, that is now superseded by finite counterexamples; but had the finite counterexamples been discovered first, we would not have found the theoretical arguments. It may be productive for future work to revisit some portions of the implication graph and locate alternate proofs and methods.

ACKNOWLEDGMENTS

Thanks to Claudio Moroni for the exploration of directed link prediction on the implication graph using GNN autoencoders described in Section 10.1.

We are also grateful to the many additional participants to the Equational Theories Project for their numerous comments and encouragement, including but not restricted to Edward van de Meent, Stanley Burris, ...

APPENDIX A. NUMBERING SYSTEM

In this section we record the numbering conventions we use for equational laws.

For this formal definition we use the natural numbers $0, 1, 2, \dots$ to represent and order indeterminate variables; however, in the main text, we use the symbols $x, y, z, w, u, v, r, s, t$ instead (and do not consider any laws with more than eight variables).

To define the ordering we use on equational laws, we first consider the case where there is a single indeterminate $*$. We place a well-ordering on words w, w' with a single indeterminate $*$ by declaring $w > w'$ if one of the following holds:

- w has a larger order than w' .
- $w = w_1 \diamond w_2$ and $w' = w'_1 \diamond w'_2$ have the same order $n \geq 1$ with $w_1 > w'_1$.
- $w = w_1 \diamond w_2$ and $w' = w'_1 \diamond w'_2$ have the same order $n \geq 1$ with $w_1 = w'_1$ and $w_2 > w'_2$.

Thus, for instance

$$* < * \diamond * < * \diamond (* \diamond *) < (* \diamond *) \diamond *.$$

We similarly place a well-ordering on equational laws $w_1 \simeq w_2$ with a single indeterminate $*$ by declaring $w_1 \simeq w_2 > w'_1 \simeq w'_2$ if one of the following holds: as follows:

- $w_1 \simeq w_2$ has a longer order than $w'_1 \simeq w'_2$.
- If $w_1 \simeq w_2$ has the same order as $w'_1 \simeq w'_2$, and $w_1 > w'_1$.
- If $w_1 \simeq w_2$ has the same order as $w'_1 \simeq w'_2$, $w_1 = w'_1$, and $w_2 > w'_2$.

Thus for instance

$$(* \diamond * \simeq * \diamond (* \diamond *)) < (* \diamond * \simeq (* \diamond *) \diamond *).$$

Finally for equational laws with alphabet $x, y, z, w, u, v, r, s, t$, define the *shape* of that law to be the law formed by replacing all indeterminates with $*$; for instance, the shape of (E4512), $(* \diamond *) \diamond * \simeq * \diamond (* \diamond *)$, is $(* \diamond *) \diamond * \simeq * \diamond (* \diamond *)$. We then place a well-ordering $w_1 \simeq w_2$ with indeterminates $x, y, z, w, u, v, r, s, t$ by declaring $w_1 \simeq w_2 > w'_1 \simeq w'_2$ if one of the following holds:

- The shape of $w_1 \simeq w_2$ is greater than the shape of $w'_1 \simeq w'_2$.
- $w_1 \simeq w_2$ and $w'_1 \simeq w'_2$ have the same shape, and the string of variables appearing in $w_1 \simeq w_2$ is lower in the lexicographical ordering (using $x < y < z < w < u < v < r < s < t$) than the corresponding string for $w'_1 \simeq w'_2$.

Thus for instance any law of shape $* \diamond * \simeq * \diamond (* \diamond *)$ is lower than any law of shape $* \diamond * \simeq (* \diamond *) \diamond *$. Among the laws of shape $* \diamond * \simeq * \diamond (* \diamond *)$, the lowest is $x \diamond x \simeq x \diamond (x \diamond x)$, which is less than (say) $x \diamond x \simeq y \diamond (y \diamond y)$, which is in turn less than $x \diamond y \simeq x \diamond (x \diamond x)$.

Two equational laws are equivalent if one can be obtained from another by some combination of relabeling the variables and applying the symmetric law $w_1 \simeq w_2 \iff w_2 \simeq w_1$. For instance, $(0 \diamond 1) \diamond 2 \simeq 1$ is equivalent to $0 \simeq (1 \diamond 0) \diamond 2$. We then replace every equational law with their minimal element in their equivalence class, which can be viewed as the *normal form* for that law; for instance, the normal form of $(0 \diamond 1) \diamond 2 \simeq 1$ would be $0 \simeq (1 \diamond 0) \diamond 2$. Finally, we eliminate any law of the form $w \simeq w$ other than $0 \simeq 0$. We then number the remaining equations $E1, E2, \dots$. For instance, $E1$ is the trivial law $0 \simeq 0$, $E2$ is the constant law $0 \simeq 1$, $E3$ is the idempotent law $0 \simeq 0 \diamond 0$, and so forth. Lists and code for generating these equations, or the equation number attached to a given equation, can be found in the ETP repository.

The number of equations in this list of order $n = 0, 1, 2, \dots$ is given by

$$2, 5, 39, 364, 4284, 57882, 888365, \dots$$

(<https://oeis.org/A376640>). The number can be computed to be

$$C_{n+1}B_{n+2}/2$$

if n is odd, 2 if $n = 0$, and

$$(C_{n+1}B_{n+2} + C_{n/2}(2D_{n+2} - B_{n+2}))/2 - C_{n/2}B_{n/2+1}$$

if $n > 2$ is even, where C_n, B_n are the Catalan and Bell numbers, and D_n is the number of partitions of $[n]$ up to reflection, which for $n = 0, 1, 2, \dots$ is

$$1, 1, 2, 4, 11, 32, 117, \dots$$

(<https://oeis.org/A103293>). A proof of this claim can be found in the ETP blueprint. In particular, there are 4694 equations of order at most 4.

Below we record some specific equations appearing in this paper, using the alphabet x, y, z, w, u, v in place of $0, 1, 2, 3, 4, 5, \dots$ for readability.

| | | |
|---------|--|--------------------------|
| (E1) | $x \simeq x$ | (Trivial law) |
| (E2) | $x \simeq y$ | (Singleton law) |
| (E3) | $x \simeq x \diamond x$ | (Idempotent law) |
| (E4) | $x \simeq x \diamond y$ | (Left-absorptive law) |
| (E5) | $x \simeq y \diamond x$ | (Right-absorptive law) |
| (E10) | $x \simeq x \diamond (y \diamond x)$ | |
| (E23) | $x \simeq (x \diamond x) \diamond x$ | |
| (E41) | $x \diamond x \simeq y \diamond z$ | |
| (E43) | $x \diamond y \simeq y \diamond x$ | (Commutative law) |
| (E46) | $x \diamond y \simeq z \diamond w$ | (Constant law) |
| (E47) | $x \simeq x \diamond (x \diamond (x \diamond x))$ | |
| (E73) | $x \simeq y \diamond (y \diamond (x \diamond y))$ | |
| (E151) | $x \simeq (x \diamond x) \diamond (x \diamond x)$ | |
| (E168) | $x \simeq (y \diamond x) \diamond (x \diamond z)$ | (Central groupoid law) |
| (E206) | $x \simeq (x \diamond (x \diamond y)) \diamond y$ | |
| (E255) | $x \simeq ((x \diamond x) \diamond x) \diamond x$ | |
| (E327) | $x \diamond y \simeq x \diamond (y \diamond z)$ | |
| (E378) | $x \simeq (x \diamond y) \diamond y$ | |
| (E395) | $x \diamond y \simeq (z \diamond x) \diamond y$ | |
| (E543) | $x \simeq y \diamond ((z \diamond (x \diamond (y \diamond z))))$ | (Tarski's axiom) |
| (E677) | $x \simeq y \diamond (x \diamond ((y \diamond x) \diamond y))$ | |
| (E817) | $x \simeq x \diamond ((x \diamond x) \diamond (x \diamond x))$ | |
| (E854) | $x \simeq x \diamond ((y \diamond z) \diamond (x \diamond z))$ | |
| (E1110) | $x \simeq y \diamond ((y \diamond (x \diamond x)) \diamond y)$ | |
| (E1117) | $x \simeq y \diamond ((y \diamond (x \diamond z)) \diamond z)$ | |
| (E1286) | $x \simeq y \diamond (((x \diamond y) \diamond x) \diamond y)$ | |
| (E1485) | $x \simeq (y \diamond x) \diamond (x \diamond (z \diamond y))$ | (Weak central groupoids) |
| (E1571) | $x \simeq (y \diamond z) \diamond (y \diamond (x \diamond z))$ | (Exp. 2 abelian groups) |
| (E1629) | $x \simeq (x \diamond x) \diamond ((x \diamond x) \diamond x)$ | |
| (E1648) | $x \simeq (x \diamond y) \diamond ((x \diamond y) \diamond y)$ | |
| (E1659) | $x \simeq (x \diamond y) \diamond ((y \diamond y) \diamond z)$ | |
| (E1689) | $x \simeq (y \diamond x) \diamond ((x \diamond z) \diamond z)$ | |
| (E1729) | $x \simeq (y \diamond y) \diamond ((y \diamond x) \diamond y)$ | |
| (E2301) | $x \simeq (y \diamond (x \diamond (y \diamond x))) \diamond y$ | |

| | | |
|-------------|---|----------------------|
| (E2441) | $x \simeq (x \diamond ((x \diamond x) \diamond x)) \diamond x$ | |
| (E2910) | $x \simeq ((y \diamond (x \diamond y)) \diamond x) \diamond y$ | |
| (E3316) | $x \diamond y \simeq x \diamond (y \diamond (x \diamond y))$ | |
| (E4315) | $x \diamond (y \diamond x) \simeq x \diamond (y \diamond z)$ | |
| (E4380) | $x \diamond (x \diamond x) \simeq (x \diamond x) \diamond x$ | |
| (E4482) | $x \diamond (y \diamond y) = (y \diamond y) \diamond x$ | |
| (E4512) | $(x \diamond y) \diamond z \simeq x \diamond (y \diamond z)$ | (Associative law) |
| (E4531) | $x \diamond (y \diamond z) \simeq (y \diamond z) \diamond x$ | |
| (E5093) | $x \simeq y \diamond ((y \diamond (y \diamond x)) \diamond (z \diamond y))$ | |
| (E345169) | $x \simeq (y \diamond ((x \diamond y) \diamond y)) \diamond (x \diamond (z \diamond y))$ | (Sheffer stroke) |
| (E42323216) | $x \simeq y \diamond (((y \diamond y) \diamond x) \diamond z)$ $\diamond (((y \diamond y) \diamond y) \diamond z)$ | (Division in groups) |

APPENDIX B. AUTHOR CONTRIBUTIONS

In a companion document to this paper, the contributions of each author of this paper to the ETP are described, following the standard CRediT categories²⁸. Below are the affiliations and grant acknowledgments of individual participants.

- Matthew Bolan: University of Toronto, matthew.bolan@mail.utoronto.ca
- Joachim Breitner: ...
- Jose Brox: ...
- Mario Carneiro: ...
- Martin Dvořák: Institute of Science and Technology Austria, martin.dvorak@matfyz.cz
- Andrés Goens: University of Amsterdam, a.goens@uva.nl
- Aaron Hill: ...
- Harald Husum: harald.husum@gmail.com
- Zoltan A. Kocsis: University of New South Wales, z.kocsis@unsw.edu.au
- Bruno Le Floch: CNRS and Laboratoire de Physique Théorique et Hautes Énergies, Sorbonne Université, blefloch@lpthe.jussieu.fr
- Lorenzo Luccioli: University of Bologna, lorenzo.luccioli2@unibo.it
- Alex Meiburg: ...
- Pietro Monticone: University of Trento, pietro.monticone@studenti.unitn.it
- Giovanni Paolini: University of Bologna, g.paolini@unibo.it
- Bernhard Reinke: ...
- David Renshaw: ...
- Marcus Rossel: Barkhausen Institut, marcus.rossel@barkhauseninstitut.org
- Cody Roux: ...
- Jérémy Scanvic, Laboratoire de Physique, École Normale Supérieure de Lyon, jeremy.scanvic@ens-lyon.fr
- Shreyas Srinivas: ...

²⁸<https://credit.niso.org/>

- Anand Rao Tadipatri: ...
- Terence Tao: Department of Mathematics, UCLA, tao@math.ucla.edu
- Vlad Tsyrklevich: vlad@tsyrklevi.ch
- Daniel Weber: ...
- Fan Zheng: ...

REFERENCES

- [1] A. K. Austin. A note on models of identities. *Proc. Amer. Math. Soc.*, 16:522–523, 1965.
- [2] A. K. Austin. Finite models for laws in two variables. *Proc. Amer. Math. Soc.*, 17:1410–1412, 1966.
- [3] Franz Baader and Tobias Nipkow. *Term rewriting and all that*. Cambridge University Press, Cambridge, 1998.
- [4] Alexander Bentkamp, Jasmin Blanchette, Visa Nummelin, Sophie Touret, Petar Vukmirovic, and Uwe Waldmann. Mechanical mathematicians. *Commun. ACM*, 66(4):80–90, 2023.
- [5] George M Bergman. The diamond lemma for ring theory. *Advances in Mathematics*, 29(2):178–218, 1978.
- [6] Joshua Clune, Yicheng Qian, Alexander Bentkamp, and Jeremy Avigad. Duper: A proof-producing superposition theorem prover for dependent type theory. In Yves Bertot, Temur Kutsia, and Michael Norrish, editors, *15th International Conference on Interactive Theorem Proving, ITP 2024, September 9-14, 2024, Tbilisi, Georgia*, volume 309 of *LIPICs*, pages 10:1–10:20. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2024.
- [7] Leonardo de Moura and Nikolaj Bjørner. Satisfiability modulo theories: An appetizer. In Marcel Vinićius Medeiros Oliveira and Jim Woodcock, editors, *Formal Methods: Foundations and Applications*, pages 23–36, Berlin, Heidelberg, 2009. Springer Berlin Heidelberg.
- [8] Leonardo Mendonça de Moura and Nikolaj S. Bjørner. Efficient e-matching for SMT solvers. In Frank Pfenning, editor, *Automated Deduction - CADE-21, 21st International Conference on Automated Deduction, Bremen, Germany, July 17-20, 2007, Proceedings*, volume 4603 of *Lecture Notes in Computer Science*, pages 183–198. Springer, 2007.
- [9] Leonardo Mendonça de Moura and Nikolaj S. Bjørner. Z3: an efficient SMT solver. In C. R. Ramakrishnan and Jakob Rehof, editors, *Tools and Algorithms for the Construction and Analysis of Systems, 14th International Conference, TACAS 2008, Held as Part of the Joint European Conferences on Theory and Practice of Software, ETAPS 2008, Budapest, Hungary, March 29-April 6, 2008. Proceedings*, volume 4963 of *Lecture Notes in Computer Science*, pages 337–340. Springer, 2008.
- [10] Gabriel Ebner, Sebastian Ullrich, Jared Roesch, Jeremy Avigad, and Leonardo de Moura. A metaprogramming framework for formal verification. *Proc. ACM Program. Lang.*, 1(ICFP):34:1–34:29, 2017.
- [11] Trevor Evans. Products of points—some simple algebras and their identities. *Amer. Math. Monthly*, 74:362–372, 1967.
- [12] Yifan Feng, Haoxuan You, Zizhao Zhang, Rongrong Ji, and Yue Gao. Hypergraph neural networks. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(01):3558–3565, July 2019.
- [13] Timothy Gowers and Michael Nielsen. Massively collaborative mathematics. *Nature*, 461(7266):879–881, October 2009.
- [14] Graham Higman and B. H. Neumann. Groups as groupoids with one law. *Publ. Math. Debrecen*, 2:215–221, 1952.
- [15] Thomas N. Kipf and Max Welling. Variational graph auto-encoders. 2016.
- [16] A. Kisielewicz. Austin identities. *Algebra Universalis*, 38(3):324–328, 1997.
- [17] Andrzej Kisielewicz. Varieties of algebras with no nontrivial finite members. In *Lattices, semigroups, and universal algebra (Lisbon, 1988)*, pages 129–136. Plenum, New York, 1990.
- [18] M. Kivela, A. Arenas, M. Barthélemy, J. P. Gleeson, Y. Moreno, and M. A. Porter. Multilayer networks. *Journal of Complex Networks*, 2(3):203–271, July 2014.
- [19] Donald E. Knuth. Notes on central groupoids. *J. Combinatorial Theory*, 8:376–390, 1970.

- [20] Donald E. Knuth and Peter B. Bendix. Simple word problems in universal algebras. In *Computational Problems in Abstract Algebra (Proc. Conf., Oxford, 1967)*, pages 263–297. Pergamon, Oxford-New York-Toronto, Ont., 1970.
- [21] Thomas Koehler, Andrés Goens, Siddharth Bhat, Tobias Grosser, Phil Trinder, and Michel Steuwer. Guided equality saturation. *Proc. ACM Program. Lang.*, 8(POPL):1727–1758, 2024.
- [22] Georgios Kollias, Vasileios Kalantzis, Tsuyoshi Idé, Aurélie Lozano, and Naoki Abe. Directed graph auto-encoders. 2022.
- [23] Laura Kovács and Andrei Voronkov. First-order theorem proving and vampire. In Natasha Sharygina and Helmut Veith, editors, *Computer Aided Verification - 25th International Conference, CAV 2013, Saint Petersburg, Russia, July 13-19, 2013. Proceedings*, volume 8044 of *Lecture Notes in Computer Science*, pages 1–35. Springer, 2013.
- [24] Daniel Kroening and Ofer Strichman. *Decision Procedures - An Algorithmic Point of View, Second Edition*. Texts in Theoretical Computer Science. An EATCS Series. Springer, 2016.
- [25] André Kündgen, Gregor Leander, and Carsten Thomassen. Switchings, extensions, and reductions in central digraphs. *J. Combin. Theory Ser. A*, 118(7):2025–2034, 2011.
- [26] Juanhui Li, Harry Shomer, Haitao Mao, Shenglai Zeng, Yao Ma, Neil Shah, Jiliang Tang, and Dawei Yin. Evaluating graph neural networks for link prediction: Current pitfalls and new benchmarking. 2023.
- [27] Jannis Limperg and Asta Halkjær From. Aesop: White-box best-first proof search for lean. In Robbert Krebbers, Dmitriy Traytel, Brigitte Pientka, and Steve Zdancewic, editors, *Proceedings of the 12th ACM SIGPLAN International Conference on Certified Programs and Proofs, CPP 2023, Boston, MA, USA, January 16-17, 2023*, pages 253–266. ACM, 2023.
- [28] W. McCune. Prover9 and mace4. <http://www.cs.unm.edu/~mccune/prover9/>, 2005–2010.
- [29] William McCune. Solution of the Robbins problem. *J. Automat. Reason.*, 19(3):263–276, 1997.
- [30] William McCune. Single axioms: with and without computers. In *Computer mathematics (Chiang Mai, 2000)*, volume 8 of *Lecture Notes Ser. Comput.*, pages 83–89. World Sci. Publ., River Edge, NJ, 2000.
- [31] William McCune, Robert Veroff, Branden Fitelson, Kenneth Harris, Andrew Feist, and Larry Wos. Short single axioms for Boolean algebra. *J. Automat. Reason.*, 29(1):1–16, 2002.
- [32] Ralph McKenzie. On spectra, and the negative solution of the decision problem for identities having a finite nontrivial model. *J. Symbolic Logic*, 40:186–196, 1975.
- [33] N. S. Mendelsohn and R. Padmanabhan. Minimal identities for Boolean groups. *J. Algebra*, 34:451–457, 1975.
- [34] C. A. Meredith and A. N. Prior. Equational logic. *Notre Dame J. Formal Logic*, 9:212–226, 1968.
- [35] V. L. Murskiĭ. The existence in the three-valued logic of a closed class with a finite basis having no finite complete system of identities. *Dokl. Akad. Nauk SSSR*, 163:815–818, 1965.
- [36] V. L. Murskiĭ. The existence of a finite basis of identities, and other properties of “almost all” finite algebras. *Problemy Kibernet.*, (30):43–56, 1975.
- [37] J. D. Phillips and Petr Vojtěchovský. The varieties of loops of Bol-Moufang type. *Algebra Universalis*, 54(3):259–271, 2005.
- [38] Marcus Rossel. An equality saturation tactic for lean. 2024.
- [39] Marcus Rossel and Andrés Goens. Bridging syntax and semantics of lean expressions in e-graphs. *arXiv preprint arXiv:2405.10188*, 2024.
- [40] Guillaume Salha, Stratis Limnios, Romain Hennequin, Viet Anh Tran, and Michalis Vazirgiannis. Gravity-inspired graph autoencoders for directed link prediction. 2019.
- [41] Th. Skolem. Logisch-kombinatorische Untersuchungen über die Erfüllbarkeit oder Beweisbarkeit mathematischer Sätze nebst einem Theoreme über dichte Mengen. *Krist. Vid. Selsk. Skr. I*, 1920, Nr. 4, 36 S. (1922)., 1922.
- [42] Alfred Tarski. Ein Beitrag zur axiomatik der abelschen gruppen. *Fundamenta Mathematicae*, 30(1):253–256, 1938.
- [43] The Coq Development Team. The Coq Proof Assistant 8.20.0, December 2024. <https://doi.org/10.5281/zenodo.14542673>.
- [44] The bbchallenge Collaboration, Justin Blanchard, Dan Briggs, Konrad Deka, Nathan Fenner, Yannick Forster, Georgi Georgiev (Skelet), Matthew L. House, Rachel Hunter, Iijil, Maja Kaździołka, Pavel Kropitz, Shawn Ligocki, mx dys, Mateusz Naściszewski, savask, Tristan Stérin, Chris Xu, Jason Yuen,

- and Théo Zimmermann. Determination of the fifth Busy Beaver value, 2025. In preparation, <https://github.com/bbchallenge/bbchallenge-paper>.
- [45] Eric Wieser. Multiple-inheritance hazards in dependently-typed algebraic hierarchies. In Catherine Dubois and Manfred Kerber, editors, *Intelligent Computer Mathematics - 16th International Conference, CICM 2023, Cambridge, UK, September 5-8, 2023, Proceedings*, volume 14101 of *Lecture Notes in Computer Science*, pages 222–236. Springer, 2023.
- [46] Max Willsey, Chandrakana Nandi, Yisu Remy Wang, Oliver Flatt, Zachary Tatlock, and Pavel Panchekha. egg: Fast and extensible equality saturation. *Proc. ACM Program. Lang.*, 5(POPL):1–29, 2021.
- [47] Stephen Wolfram. The physicalization of metamathematics and its implications for the foundations of mathematics, Mar 2022.
- [48] Michihiro Yasunaga, Antoine Bosselut, Hongyu Ren, Xikun Zhang, Christopher D Manning, Percy Liang, and Jure Leskovec. Deep bidirectional language-knowledge graph pretraining. 2022.
- [49] Xikun Zhang, Antoine Bosselut, Michihiro Yasunaga, Hongyu Ren, Percy Liang, Christopher D. Manning, and Jure Leskovec. Greaselm: Graph reasoning enhanced language models for question answering. 2022.
- [50] Xitong Zhang, Yixuan He, Nathan Brugnone, Michael Perlmutter, and Matthew Hirn. Magnet: A neural network for directed graphs. 2021.

TODO LIST

| | |
|---|----|
| This is my vague understanding of the situation - someone with experience should validate this. | 13 |
| Shreyas Srinivas and Pietro Monticone have volunteered to take the lead on this section. | 13 |
| Discuss how the above breaks down when a large and diverse collection of collaborators work over the internet. Consider citing the EMS webpage for code of practice for joint responsibility rules. Start the project organisation setup in the next subsection | 13 |
| elaborate on what level of overlap was permissible and what we consider excessive . . . | 15 |
| Discuss this further, perhaps give an example, or refer to the ATP section. See also the discussion threads “Counterexamples by enumerating words in quotient magmas” and “Using SAT solvers for model generation” threads (sorry, LaTeX is choking on the URLs for some reason). | 16 |
| Discuss performance of this method. | 18 |
| Give some statistics of what proportion of refutations can be resolved by linear models. | 19 |
| Report on how large the twisting semigroups are in practice, and how many implications can be refuted by this method. | 22 |

| | |
|--|----|
| Describe performance of this automated method. Discuss the issue that some implications required a large SAT solver calculation that was difficult to formalize efficiently in Lean, prompting human-generated simplified proofs using smaller rulesets. | 24 |
| here we could add a screenshot of using vampire or prover9 | 37 |
| Improve the structure of this section | 38 |
| add explanations from the following URL | 38 |
| this should also be expanded more, at least with references to some of the constructions in other chapters. | 40 |
| discuss guided search in the finite implications or the Higman-Neumann work Jose Brox has done. | 40 |
| maybe add a screenshot here of the workflow of using a seed to find counterexamples with prover9 or vampire? | 40 |
| throw a couple of "benchmarking" tables for the same ATP with different parameters and for different ATPs, talk about some relative gains in time (changing parameters we saw a 500 times speedup on this particular problem), etc. This is knowledge I think we have gained to some extent, and certainly I would have been glad to receive this kind of hints before we started!". Then leave it as an interesting open problem to properly develop and measure benchmarks for ATPs based on this project. | 40 |
| Expand this sketch | 41 |
| report on Higman-Neumann laws | 41 |
| Describe how data was handled during the project and how it will be managed going forward. | 48 |